



Proceedings of the 25th Annual Conference of the European Association for Machine Translation

Volume 1: Research And Implementations & Case Studies

June 24-27, 2024
Sheffield, United Kingdom

Edited by:

Carolina Scarton, Charlotte Prescott, Chris Bayliss, Chris Oakley, Joanna Wright, Stuart Wrigley, Xingyi Song, Edward Gow-Smith, Rachel Bawden, Víctor M Sánchez-Cartagena, Patrick Cadwell, Ekaterina Lapshinova-Koltunski, Vera Cabarrão, Konstantinos Chatzitheodorou, Mary Nurminen, Diptesh Kanojia, Helena Moniz

Organised by





The papers published in this proceedings are —unless indicated otherwise— covered by the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International (CC-BY-NCND 4.0). You may copy, distribute, and transmit the work, provided that you attribute it (authorship, proceedings, publisher) in the manner specified by the author(s) or licensor(s), and that you do not use it for commercial purposes. The full text of the licence may be found at <https://creativecommons.org/licenses/by-nc-nd/4.0/deed.en>

©2024 The authors

ISBN 978-1-0686907-0-9

Publisher: European Association for Machine Translation (EAMT)

Foreword from the General Chair

As president of the European Association for Machine Translation (EAMT) and General Chair of the 25th Annual Conference of the EAMT, it is with great pleasure that I write these opening words to the Proceedings of EAMT 2024, a special year since we are celebrating our 25th anniversary!

According to tradition, my first note of deep appreciation and gratitude goes to Celia Rico, Luc Meertens, Lucia Specia, and Maja Popovič, Executive Board Members, who have moved to new adventures in their lives, after outstanding, and dedicated service to the EAMT community.

We have several milestones to celebrate this year, built upon the hard work of our Executive Committee (EC) and our community: upgraded grants for low-income and war zones and for Translation Studies, a record submission rate for research projects, continuous excellent submissions for the best thesis award, and one of the highest number of papers ever submitted to our conference (80 papers accepted)! I could not be prouder of our EC and the dynamics of our community.

The EAMT Executive Committee (EC) has been very busy. Luc Meertens (treasurer), Carolina Scarton (secretary) and Sara Szoc (preparing to become our secretary and supporting everything we do) have been tirelessly supporting all initiatives. André Martins and Celia Rico, our co-chairs for low-income areas, war zones and Translation Studies grants, selected 11 grantees, 6 applicants from Translation Studies and 5 from war zones (3 hybrid light and 8 in-person). Maja Popovič and Sara Szoc, our co-chairs for the Research Projects, selected 4 projects (equally distributed by students and general research projects calls) with a diverse set of topics. To all our co-chairs, my gratitude! The selection work is never an easy task and this year was particularly hard.

The same applied to the best thesis award – Barry Haddow, chair of the Best Thesis Award, had a very difficult time selecting a candidate, since the submissions were of very high quality. Our congratulations to Marco Gaido's thesis "Direct Speech Translation Toward High-Quality, Inclusive, and Augmented Systems" (FBK, Italy), supervised by Marco Turchi and Matteo Negri. Our congratulations extended to the two highly commended theses of Jannis Vamvas: "Model-based Evaluation of Multilinguality" (University of Zurich, Switzerland), supervised by Rico Sennrich and Lena A. Jäger; and Javier Iranzo-Sánchez: "Streaming Neural Speech Translation" (UPV, Spain), supervised by Jorge Civera and Alfons Juan.

EAMT, as full sponsor of the MT Marathon, would also like to highlight the outstanding work that the MT Marathon organisers conducted, enriching the vitality of our community with their projects and keynotes. A special thank you to the organising committee Lisa Yankovskaya, Agnes Luhtaru, Lisa Korotkova, Mark Fišel, Ondrej Bojar, and Barry Haddow for all the efforts on yet another successful MT Marathon event. Thank you, University of Tartu, for hosting the event.

Sheffield, United Kingdom! EAMT 2024 celebrates our 25th anniversary! Our conference will have a three-day, four-track programme put together by our chairs: Rachel Bawden and Víctor Sánchez-Cartagena (research: technical track co-chairs); Ekaterina Lapshinova-Koltunski and Patrick Cadwell (research: translators & users track co-chairs); Chatzitheodorou Konstantinos and Vera Cabarrão (implementations & case studies track co-chairs); and Mikel Forcada and Helena Moniz (products & projects track chairs). And backing up all the scientific components of our conference and filters of quality for the final selection: our reviewers. Thank you for your work and the alignment between all the chairs!

Continuing the successful event from Tampere, this year EAMT 2024 will also have an extra day for workshops and tutorials, organised by our co-chairs Diptesh Kanojia and Mary Nurminen. Once more, the submissions for workshops and tutorials largely exceeded our expectations for our second edition!

The programme will continue the tradition of including two keynote speakers, Alexandra Birch (Reader in Natural Language Processing in the Institute for Language, Cognition and Computation, School of Informatics, University of Edinburgh) and Valter Mavrič (Director-General of the Translation Service – DG TRAD – at the European Parliament). Our outstanding keynote speakers will demonstrate their extensive and global impactful work in translation studies and translation technologies.

EAMT 2024 would never be possible without the synergetic, sharp, enthusiastic, and hard working local organising team! What a dream and fun team to work with! Our local co-chair, Carolina Scarton (University of Sheffield, UK), who always supports the EAMT community and is always eager to do the best EAMT ever! Our local co-chair from ZOO Digital, Chris Oakley, also Charlotte Prescott (ZOO Digital, UK), Chris Bayliss (ZOO Digital, UK), Joanna Wright (University of Sheffield, UK), and Xingyi Song (University of Sheffield, UK). From the local organising support team, our thank you to Freddy Heppell (University of Sheffield, UK) and Tom Pickard (University of Sheffield, UK). Our special gratitude to the University of Sheffield and ZOO Digital for the joint efforts. You will surely make our 25th anniversary memorable!

The Sheffield team is working towards a special 25th anniversary. Carolina Scarton has been doing intensive work on organising and finding a home for the John Hutchins Machine Translation Archive. Carolina is deeply committed to respect John's wishes of making his library available to the community, and the former president, Mikel Forcada, and current one are fully supporting Carolina's initiatives. As an anticipation of such effort, the Sheffield team is working on presenting a sample of John's books for EAMT 2024 participants! Thank you, Carolina Scarton, for all the hard work on this. Within this topic still, a special thank you to Mike Hutchins, John's son, who is fully committed to make it happen and respect his father's vision of giving back to the community.

EAMT has been supported by generous sponsors in its initiatives along the years. This year is no exception. Our gratitude to our Silver sponsors: RWS Language Weaver, Translated, and Unbabel. To our Bronze sponsors: CrossLang, Pangeanic, STAR, and TransPerfect. Also to Apertium, our long standing collaborator sponsor, Springer, our Supporter sponsor for the Best Paper award, and our Media sponsors, MultiLingual. Your support is vital in our efforts to give back to our community through grants and other initiatives.

A note still to all our EAMT members and our participants! Without you no effort would make sense! Let us take this opportunity to create scientific collaboration and give constructive feedback. To fully enjoy the conference, please check our Code of Conduct at <https://eamt2024.sheffield.ac.uk/code-of-conduct>. I'm looking forward to seeing you all and celebrating our 25th anniversary with you!

It is our organisation's greatest wish to continue giving back to our community and to drive and be driven by our community's energy and enthusiasm. Reach out to us if you have new ideas or suggestions you would like to implement. We will try hard to accomplish it with you. Learn more about us at <https://eamt.org/>.

Helena Moniz

President of the EAMT
General Chair of EAMT 2024
University of Lisbon / INESC-ID, Portugal

Message from the Organising Committee

Ey Up!

We are delighted to welcome you to EAMT 2024 at Sheffield and celebrate its 25th anniversary. Sheffield, renowned for its rich industrial heritage and pivotal role in the steel industry, provides an ideal venue for “forging” collaboration and exchanging ideas. The outdoor city provides an ideal and welcoming environment for a thriving international community with a large number of students. The UK’s greenest city has the Peak District National Park at its doorstep, being a not to be missed place for the most adventurous (looking for sports like bouldering and mountain biking) as well as for just relaxing on a short walk enjoying the views and hospitality of the Peak District’s small villages. It is not rare that students end up staying in Sheffield and calling this fabulous place home (which is the case of some of us on the organising committee).

The University of Sheffield has also been key in developing Machine Translation research, being an active member of EAMT and part of its history. Memorable former members of the Sheffield community include: the late John Hutchins (creator of the MT Archive and author of the 1992 book *An introduction to machine translation*) was a librarian in Sheffield from 1965 and 1971; the late Professor Yorick Wilks (author of the 2008 book *Machine Translation: Its Scope and Limits*) was an emeritus professor and a former Head of the Computer Science department; and Professor Lucia Specia (the pioneer in the area of MT Quality Estimation and author of the 2018 book *Quality Estimation for Machine Translation*) was professor at the Computer Science department and former PhD supervisor of two of the local organisers.

ZOO Digital is a global provider of cloud-based localisation and digital distribution services for the media and entertainment industry. ZOO Digital offers a range of services including subtitling, dubbing, media processing, and distribution. The company uses proprietary technology platforms to streamline and manage the localisation process, making it more efficient and cost-effective. ZOO is a long-term partner of the University of Sheffield, being committed to support research in speech and text translation. They are also one of the most active sponsors of our UKRI AI Centre for Doctoral Training (CDT) in Speech and Language Technologies and their Applications and had their first sponsored PhD student working on the area of MT graduating in 2023.

We are especially excited about our conference venues, which showcase some of Sheffield’s most iconic sites. Our welcome reception will take place in the stunning Sheffield Winter Garden, one of the largest temperate glasshouses in the UK. This beautiful indoor garden is filled with exotic plants from around the world. The conference dinner will be hosted at the Kelham Island Museum, a celebrated institution that chronicles the city’s industrial history and innovation in steel production. Attendees will have the unique opportunity to visit the impressive River Don Engine, a steam engine that highlights Sheffield’s engineering and industrial heritage. We are also thrilled to announce that ZOO Digital has generously funded a special pre-conference social event at the National Videogame Museum. This interactive museum celebrates the history and culture of video games, offering a fun and engaging way for attendees to unwind and connect with each other. Finally, participants that opt to attend the Kelham Island Food tour will be taken on a culinary journey of the area, visiting a range of eating establishments and enjoying generous samples at each stop, and gaining insight into the interesting history of this famous Sheffield district.

We extend our deepest gratitude to our Silver Sponsors (Language Weaver, Translated, Unbabel), Bronze Sponsors (AppTek, CrossLang, Pangeanic, STAR Group, TransPerfect), Collaborator (Apertium), Sponsor (Springer Nature), Media Sponsors (MultiLingual), track chairs (Helena Moniz, Rachel Bawden, Víctor M Sánchez-Cartagena, Patrick Cadwell, Ekaterina Lapshinova-Koltunski, Vera Cabarrão, Konstantinos Chatzitheodorou, Mikel Forcada, Mary Nurminen, Diptesh Kanojia, Barry Haddow), keynote speakers (Alexandra Birch, Valter Mavrič), the programme committee, and authors.

Our special very thanks goes to the volunteers (Freddy Heppell, Tom Pickard, Edward Gow-Smith, and Shenbin Qian), administrative and technical support (Natalie Hothersall, Kim Matthews-Hyde, and James Bishop), events management (Gavin Lambert), and our emergency organisation support committee (Xi Wang and Mark Stevenson) whose hard work and dedication have made this conference possible. We also thank the EAMT executive committee for all the support provided and trust in our work, in particular Helena Moniz (also our general chair) and Sara Szoc. Finally, we also thank the Department of Computer Science, in particular Professor Heidi Christensen (Head of the Computer Science department) and Professor Kalina Bontcheva (head of the Natural Language Processing research group), for their support of our conference.

We invite you to explore and enjoy the city of Sheffield. Whether you are discovering its historical landmarks, enjoying its green spaces, or immersing yourself in its rich cultural offerings, we hope you find inspiration both within and beyond the conference sessions.

Carolina Scarton
(University of Sheffield)
(EAMT Secretary)

Charlotte Prescott
(ZOO Digital)

Chris Bayliss
(ZOO Digital)

Chris Oakley
(ZOO Digital)

Joanna Wright
(University of Sheffield)

Stuart Wrigley
(University of Sheffield)

Xingyi Song
(University of Sheffield)

Preface by the Programme Chairs

On behalf of the programme chairs, a warm welcome to the 25th annual conference of the European Association for Machine Translation in Sheffield, UK. Following last year's restructuring of the research track into two tracks, this year's conference programme is divided into four tracks, two dedicated to research (one for technical papers for development of MT techniques and one focused on translators and users of MT), an implementations and case studies track and a projects and products track.

The **Technical Research track** invited submissions on significant results in any aspect of MT and related areas, including multilingual technologies. As in previous years, this track proved the most popular of the four tracks, receiving a total of 46 submissions from 26 different countries. With one desk rejection and four paper withdrawals, 20 papers were accepted from 18 different countries, resulting in an acceptance rate of 43%, which is consistent with previous years. Six of the accepted papers are to be presented orally and the remaining 14 will be presented as posters.

Following current practices in the field, papers focus on neural MT (NMT), with several works also studying large language models (LLMs) for translation. Accepted papers represented a wide range of topics relevant to current interests in the field: context-aware MT (Appicharla et al., 2024; Gete and Etchegoyhen, 2024); the application of techniques for low-resource languages and scenarios (Chen et al. 2024; Guttman et al.; Simonsen and Einarsson, 2024; Song et al. 2024) including sign language translation (McGill et al., 2024); attention to specific domains (Ploeger et al., 2024; Roussis et al. 2024) and to the challenges faced when dealing with them, e.g. for the incorporating of terminologies (Hauhio and Friberg. 2024). A number of works study LLMs (Chen et al., 2024.; Mujadia et al. 2024; Simonsen and Einarsson, 2024), a trend that is likely to continue in years to come. As a sign of the progress being made in the quality of MT systems, the EAMT 2024 technical research track also features several papers dealing with topics related to the alignment of MT outputs with the expectations of human users (Moura Ramos et al., 2024), including on the topics of toxicity (García Gilabert et al., 2024), formality (Wisniewski et al., 2024) and gender-inclusiveness (Piergentile et al., 2024).

We would like to give our thanks to all the authors who submitted to the track and to the 72 reviewers, who provided feedback and insightful comments for the submissions received. We are particularly grateful to the emergency reviewers who agreed to review papers at the last minute in order for decision notifications to be sent out on time.

Translators and Users Track

The focus of the Translators and Users track is to cover a wide range of topics related to the interaction between human translators and other users of machine translation. The second edition of this track attracted 21 papers, with 18 accepted out of them which comprises 85.71% of acceptance. Five of the accepted papers will be presented orally and 13 will be presented at a dedicated poster presentation session. The accepted papers address the interaction between machine translation and its users from various perspectives and cover various aspects of machine translation use, including both interlingual and intralingual translation, looking into challenges and potentials of large language models, as well as correlating human and machine translation. They provide novel examinations of long-standing areas of interest for translators and users in this space including translation quality, MT performance, tools and methods to assist translators, and users' perceptions and attitudes towards MT.

Sui He experiments with prompts applying ChatGPT for automatic translation. The author compares translation briefs and what s/he calls persona prompts (assignment of a role of an author or translator to the system).

Claudio Fantinuoli and Xiaoman Wang explore correlation between automatic quality evaluation metrics with human judgements for simultaneous interpreting.

Serge Gladkoff et al. investigate the application of the state-of-the-art LLMs for uncertainty estimation of MT output quality, which is required to determine the need for post-editing.

Paolo Canavese and Patrick Cadwell analyse translators' perspectives on the use of machine translation and its impact in a specific institutional setting, i.e. the Swiss Confederation.

Marta R. Costa-jussà et. al. presents a novel multimodal and multilingual pipeline to automatically identify and mitigate added toxicity at inference time, which does not require further model training.

Celia Soler Uguet et al. compare performance of various LLMs for automatic post-editing and MQM error annotation across four languages in a medical domain.

Lise Volkart and Pierrette Bouillon compare human translation and post-edited machine translation from a lexical and syntactic perspective in two language pairs: English-French and German-French. Their aim is to find out if NMT systems produce lexically and syntactically poorer translations.

Gabriela Gonzalez-Saez et al. describe their work on visualisation tools to foster collaborations between translators and computational scientists.

Maria Kunilovskaya et al. explore if GPT-4 can reduce translationese (specific feature of translated texts) in human-translated texts on bidirectional German-English data from the Europarl corpus.

Rachel Bawden et al. evaluate the effectiveness of a post-editing pipeline for the translation of scientific abstract demonstrating that such pipelines can be effective for high-resource language pairs.

Vicent Briva-Iglesias and Sharon O'Brien present a user study on professional English-Spanish translators in the legal domain, which focuses on impact of negative or positive translators' pre-task perceptions of MT.

Miguel Rios et al. explore the impact of automatic speech synthesis in a post-editing machine translation environment in terms of quality, productivity, and cognitive effort.

Silvana Deilen et al. evaluate performance of intralingual machine translation systems in the area of health communication.

Michael Carl looks into a way of using machine learning to validate the empirical objectivity of a taxonomy for behavioral translation data.

João Lucas Cavalheiro Camargo et al. conduct a survey aimed at identifying and exploring the attitudes and recommendations of machine translation quality assessment educators.

Bettina Hiebl and Dagmar Gromann propose to use the Best-Worst scoring for a comparative translation quality assessment of one human and three machine translations in the English-German language pair.

Adaeze Ngozi Oluoba et al. investigate methods to detect critical and harmful MT errors caused by non-compositional multi-word expressions and polysemy. For this, they design diagnostic tests that they apply on collections of medical texts.

Nora Aranberri explores evaluation of the Spanish-Basque translations. The author compares evaluations done by volunteers and translation professionals.

We would like to thank the 28 colleagues that kindly gave their time and effort to review the papers submitted to this track. Your reviews were perceptive, detailed, and, above all, constructive. We would also like to express our special gratitude to those reviewers who stepped in at the last minute to provide extra reviews at short notice. Your collegiality was a great support to us.

Implementations and case studies track

Entering the second year with the Implementations & Case Studies track, we are excited to share the acceptance of 9 papers. These papers cover a wide range of topics, showing the latest advancements, challenges, and creative ideas in MT. The goal for this track remains unchanged: to report experiences with MT in organizations of all types (both industry and academia) and to share views and observations based on day-to-day experiences working within the dynamic field of MT.

The journey begins with Oliver et al. who detail corpus creation and NMT model training for legal texts in low-resource languages, shedding light on the intricacies of bridging linguistic gaps in specialized domains.

Continuing on this path, Eschbach-Dymanus et al. delve into the realm of domain adaptation of MT for business IT texts, offering valuable insights into the translation capabilities of LLMs.

Bechara et al. present the creation and evaluation of a multilingual corpus of UN General Assembly debates, underscoring the importance of robust linguistic resources in advancing our understanding of multilingual communication.

Additionally, Korotkova and Fishel present groundbreaking research on Estonian-centric MT, emphasizing data availability and releasing a back-translation corpus of over 2 billion sentence pairs.

Moving forward, Silveira et al. examine the suitability of GPT-4 in generating subject-matter expertise assessment questions, illuminating new avenues for leveraging artificial intelligence in language assessment.

Continuing in this direction, Nunziatini et al.'s research explores the advantages and disadvantages of using LLMs to make raw MT output gender-inclusive.

Berger et al. work in prompting LLMs with human error markings represents a significant step towards self-correcting MT, offering promising avenues for enhancing translation quality in specialized domains.

Vasiljevs et al. present findings from a comprehensive market study on advancing digital language equality in Europe. They provide critical insights into the current landscape of multilingual website translation and introduce innovative open-source solutions aimed at bridging linguistic divides.

Lastly, Vincent et al. present an insightful case study on contextual MT in professional subtitling. This work sheds light on the practical implications of incorporating extra-textual context into the MT pipeline, offering valuable lessons for industry practitioners.

Together, these papers paint a vivid picture of the ever-evolving landscape of MT Implementations & Case Studies, showcasing the ingenuity, resilience, and collaborative spirit of the MT community.

Products and Projects track

This year we received 31 submissions and 30 papers were accepted. The selection will provide a plethora of products and projects being developed by our community with a rich set of topics, ranging from EAMT sponsored projects, European projects, services and products from distinguished industry and research players of our community. It will surely be a very lively session with the usual poster boosters (one of our EAMT conferences' favourite moments) and poster sessions. We would like to thank the 25 reviewers, who were drafted quite late, for their quick response and their timeliness.

Rachel Bawden
(Inria, Paris, France)

Víctor M Sánchez-Cartagena
(University of Alacant, Spain)

Patrick Cadwell
(DCU, Ireland)

Ekaterina Lapshinova-Koltunski
(University of Hildesheim, Germany)

Vera Cabarrão
(Unbabel, Portugal)

Konstantinos Chatzitheodorou
(Strategic Agenda, UK)

Helena Moniz
(University of Lisbon (FLUL)
INESC-ID, Portugal)

Mikel Forcada
(Prompsit Language Engineering
Elx, Spain)

Mary Nurminen
(Tampere University, Finland)

Diptesh Kanojia
(University of Surrey, UK)

Barry Haddow
(University of Edinburgh, UK)

EAMT 2023 Best Thesis Award (Anthony C Clarke Award)

For the 2023 best theses award, we received a total of 9 submissions; all were MT-related thesis defended in 2023. We recruited 20 reviewers to examine and score the theses, considering how challenging the problem tackled in each thesis was, how relevant the results were for machine translation as a field, and what the strength of its impact in terms of scientific publications was. Two EAMT Executive Committee members also analysed all theses. It became very clear that 2023 was another very good year for PhD theses in machine translation.

All theses had merit, all candidates had strong CVs and, therefore, it was very difficult to select a winner.

A panel of two EAMT Executive Committee members (Barry Haddow and Helena Moniz) was assembled to process the reviews and select a winner that was later ratified by the EAMT executive committee.

We are pleased to announce that the **winner of the 2023 edition of the EAMT Best Thesis Award is Marco Gaido's thesis "Direct Speech Translation Toward High-Quality, Inclusive, and Augmented Systems"** (FBK, Italy), supervised by Marco Turchi and Matteo Negri.

In addition, the committee judged that the following theses, were **"highly commended"**:

Jannis Vamvas: "Model-based Evaluation of Multilinguality" (University of Zurich, Switzerland), supervised by Rico Sennrich and Lena A. Jäger

Javier Iranzo-Sánchez: "Streaming Neural Speech Translation" (UPV, Spain), supervised by Jorge Civera and Alfons Juan

The awardee will receive a prize of €500, together with a suitably-inscribed certificate. In addition, Dr. Gaido will present a summary of their thesis at the 25th Annual Conference of the European Association for Machine Translation. In order to facilitate this, the EAMT will waive the winner's registration costs, and will make available a travel bursary of €200.

Barry Haddow, chair, EAMT BTA award 2023
University of Edinburgh, UK

Organising Committee

General Chair

Helena Moniz, Universidade de Lisboa / INESC-ID

Local Organising Committee

Carolina Scarton, University of Sheffield
Charlotte Prescott, ZOO Digital
Chris Bayliss, ZOO Digital
Chris Oakley, ZOO Digital
Joanna Wright, University of Sheffield
Xingyi Song, University of Sheffield

Local Organising Support Team

Edward GowSmith, University of Sheffield
Freddy Heppell, University of Sheffield
Tom Pickard, University of Sheffield
Shenbin Qian, University of Surrey

Implementations Case Studies Track Program Chairs

Vera Cabarrão, Unbabel
Konstantinos Chatzitheodorou, Strategic Agenda

Products and Projects Track Program Chairs

Helena Moniz, Universidade de Lisboa
Mikel Forcada, Prompsit Language Engineering

Research Translators Users Track Program Chairs

Patrick Cadwell, Dublin City University
Ekaterina Lapshinova-Koltunski, Universität Hildesheim

Technical Track Program Chairs

Rachel Bawden, Inria
V́ctor M. Sánchez-Cartagena, Universidad de Alicante

Thesis Award Program Chairs

Barry Haddow, University of Edinburgh

Workshops and Tutorials Program Chairs

Diptesh Kanojia, University of Surrey
Mary Nurminen, Tampere University

Programme Committee

Implementations Case Studies Track

Eleftherios Avramidis, Fred Bane, Adam Bittlingmayer, Marianna Buchicchio, Laura Casanellas, Laura Casanellas, Konstantin Dranch, László János Laki, Mara Nunziatini, Raj Nath Patel, Spyridon Pilos, Heather Rossi, Konstantin Savenkov, Marina Sánchez Torrón, Anna Zaretskaya

Research Translators Users Track

Sergi Alvarez-Vidal, Nora Aranberri, Lynne Bowker, Vicent Briva-Iglesias, João Lucas Cavaleiro Camargo, Michael Carl, Dragoş Ciobanu, Oliver Czulo, Joke Daems, Christophe Declercq, Dr. Silvana Deilen, Félix Do Carmo, Aletta G. Dorst, Maria Fernandez-Parra, Federico Gaspari, Junyan Jiang, Ramuné Kasperé, Dorothy Kenny, Maarit Koponen, Rudy Looock, Lieve Macken, Antoni Oliver, David Orrego-Carmona, Maria Del Mar Sánchez Ramos, Celia Rico, Carlos S C Teixeira, Susana Valdez, Mihaela Vela, Lucas Nunes Vieira

Technical Track

Sweta Agrawal, Eleftherios Avramidis, Parnia Bahar, Loic Barrault, Magdalena Biesialska, Sheila Castilho, Chloé Clavel, Éric Villemonte De La Clergerie, Raj Dabre, Aswarth Abhilash Dara, Miguel Domingo, Hiroshi Echizenya, Cristina España-Bonet, Miquel Esplà-Gomis, Marcello Federico, Marco Gaido, Aarón Galiano-Jiménez, Mattia Antonino Di Gangi, Thanh-Le Ha, Rejwanul Haque, Rebecca Knowles, Philipp Koehn, Maria Kunilovskaya, Gregor Leusch, Andreas Maletti, Antonio Valerio Miceli Barone, Kenton Murray, Masaaki Nagata, Toshiaki Nakazawa, Jan Niehues, Constantin Orasan, Daniel Ortiz-Martínez, Pavel Pecina, Stephan Peitz, Sergio Penkale, Andrei Popescu-Belis, Maja Popovic, Juan Antonio Pérez-Ortiz, Tharindu Ranasinghe, Natalia Carolina Alencar De Resende, Miguel Rios, Rudolf Rosa, Fatiha Sadat, Benoît Sagot, Beatrice Savoldi, Yves Scherrer, Djamé Seddah, Rico Sennrich, Dimitar Shterionov, Michel Simard, Patrick Simianer, Mirella De Sisto, Felix Stahlberg, Katsuhito Sudoh, Felipe Sánchez-Martínez, Aleš Tamchyna, Joël Tang, Ayla Rigouts Terryn, Arda Tezcan, Jörg Tiedemann, Antonio Toral, Masao Utiyama, Vincent Vandeghinste, Dušan Variš, David Vilar, Martin Volk, Trang Vu, Taro Watanabe, Minghao Wu, François Yvon, Biao Zhang, Dakun Zhang

Thesis Award

Rachel Bawden, Daniel Beck, Alexandra Birch, Ondřej Bojar, Bill Byrne, Vera Cabarrão, Sheila Castilho, Anna Currey, José G. C. De Souza, Miquel Esplà-Gomis, Marcello Federico, Mikel L. Forcada, Liane Guillou, Diptesh Kanojia, Philipp Koehn, Mary Nurminen, Constantin Orasan, John E. Ortega, Santanu Pal, Danielle Saunders, Carolina Scarton, Xingyi Song, Felix Stahlberg, Antonio Toral, Marina Sánchez Torrón, Bram Vanroy, Marcely Zanon Boito

Keynote Talk

Harnessing the benefits of machine translation at the European Parliament: from current practices to future possibilities

Valter Mavrič
European Parliament
24-06-2024 11:00:00

Abstract: Machine translation (MT) is an essential tool for one of the largest institutional translation providers in the world: the European Parliament's Directorate-General for Translation (DG TRAD). DG TRAD is home to 24 language units that embody and put into practice one of the core democratic principles of the European Union: multilingualism. In this complex environment, MT has become an integral part of DG TRAD's work, helping it to manage an ever-growing volume of translation requests and allowing it to focus on the unique value that only humans can bring to the translation process. The MT technology used in DG TRAD is a focal point of cooperation between the EU institutions and is constantly evolving. To best harness the benefits, DG TRAD relies on a dedicated team that carries out tests to explore the best ways of using MT for DG TRAD's content. This presentation will tell you, from a user's perspective, about DG TRAD's journey to identify the most efficient ways of working with MT. Here are some of the questions we will cover:

- How well does MT handle the European Parliament's content? Do all languages produce the same results? How does MT quality vary based on the type of content?
- How does MT improve efficiency? What efforts are still necessary after integrating MT into DG TRAD's workflow?
- What about clear language? How well does MT perform in this area?

Finally, we will look at the new areas DG TRAD is exploring in this age of artificial intelligence (AI) and where we see that further research could provide added value.

Bio: Valter Mavrič is Director-General of the Translation Service (DG TRAD) at the European Parliament (since 2016), where he was previously acting Director-General (from 2014), Director (from 2010) and Head of the Slovenian Translation Unit (from 2004). With an MA in applied linguistics and further training in translation, interpretation, linguistics and management, he has a long experience as manager, translator, interpreter and teacher of languages. He works in Slovenian, Italian, English, French, and Croatian and is currently preparing a PhD in strategic communication.

Keynote Talk

Translation and LLMs

Alexandra Birch

School of Informatics, University of Edinburgh

26-06-2024 09:15:00

Abstract: What is the future of translation research in the era of large language models? Brown et al. in 2020 showed that prompting GPT3 with a few examples of translation could result in translations which were higher quality than SOTA supervised models at the time (into English and only for French, German). Until this point, research on machine translation had been central to the field of natural language processing, often attracting the most submissions in annual NLP conferences and leading to many breakthroughs in the field. Since then, there has been enormous interest in models which can perform a wide variety of tasks and interest in translation as a separate sub-field has somewhat diminished. However, translation remain a compelling and widely used technology. So what is the promise of LLMs for translation and how should we best use them? What opportunities do LLMs unlock and what challenges remain? How can the field of translation still contribute to NLP? I will touch on some of my own research but I focus on these broader questions.

Bio: Alexandra Birch is a Reader in Natural Language Processing in the Institute for Language, Cognition and Computation (ILCC), School of Informatics, University of Edinburgh. She is a leader of the StatMT group and a co-founder of Aveni.ai - an award winning startup in speech analytics and conversational AI. Her main research focuses on machine translation and multilingual dialogue, but she has a broad interest in leveraging NLP to create compelling applications that improve people's lives.

Tutorial

Linguistically Motivated Neural Machine Translation

Haiyue Song, Hour Kaina, Raj Dabre

National Institute of Information and Communications Technology (NICT), Japan

27-06-2024 09:00:00

Abstract: In this tutorial, we focus on a niche area of neural machine translation (NMT) that aims to incorporate linguistics into different stages in the NMT pipeline, from pre-processing to model training to evaluation. We first introduce the background of NMT and fundamental analysis tools, such as word segmenters, part-of-speech taggers, and dependency parsers. We then cover topics including 1) word/subword segmentation, and character decomposition during MT data pre-processing, 2) incorporating direct and indirect linguistic features into NMT models, and 3) fine-grained linguistic evaluation for MT systems. We reveal the impact of orthography, syntax, and semantics information on translation performance. This tutorial is mainly aimed at researchers interested in the intersection of linguistics and low-resource machine translation. We hope this tutorial inspires and encourages them to develop linguistically motivated high-quality MT systems and evaluation benchmarks.

Panel

LLMs and Machine Translation for Low-Resource Languages: Bridging Gaps or Widening Divides?

24-06-2024 15:00:00 - 17:00:00

LLMs such as ChatGPT, Claude and Gemini 1.5 have come to dominate the AI landscape, through their ability to perform well across a wide range of tasks and languages. They have excellent abilities in machine translation for high-resource languages, often performing on par with dedicated translation models, and with exciting use-cases including stylization, post-editing, and human-in-the-loop approaches. Nevertheless, these models' capabilities are much more limited in languages with less digital representation: performance in lower-resource languages can be regarded as a byproduct rather than a focus and the reliance on English language training data reinforces English language cultural hegemony, with particularly high representation of American English cultural knowledge in model weights. In downstream evaluation, claims of multilinguality typically belie the dependence on English-centric data: the FLORES dataset, for example, which contains MT evaluation data in over 200 languages, is largely translated from English. This panel will explore the challenges and opportunities associated with LLMs for translating low-resource languages, investigating the dangers of exacerbating existing linguistic and cultural biases, the potential of LLMs to democratise information access, and how to ensure that these models benefit rather than marginalise underrepresented linguistic communities.

Panelists:

Adaeze Ngozi Oluoba, University of Leeds, UK Adaeze Ngozi Oluoba is a PhD researcher at the School of Languages, Cultures and Societies, University of Leeds. Her PhD research focuses on using large language models to detect and predict English medical source texts that could produce potentially harmful outputs when machine translated into a low-resource language like Igbo. Prior to commencing her PhD studies, she worked as a lecturer at the Department of Foreign Language and Translation Studies, Abia State University, Nigeria. She is also a freelance translator/ editor specialising in legal, medical and literary translations from French/Igbo into English and English/French into Igbo. Her research interests include Machine Translation for Low-Resourced Languages, Computational Linguistics, French as a Foreign Language and Language in Health

Alexandra Birch, University of Edinburgh, UK Alexandra Birch is a Reader in Natural Language Processing in the Institute for Language, Cognition and Computation (ILCC), School of Informatics, University of Edinburgh. She is a leader of the StatMT group and a co-founder of Aveni.ai - an award winning startup in speech analytics and conversational AI. Her main research focuses on machine translation and multilingual dialogue, but she has a broad interest in leveraging NLP to create compelling applications that improve people's lives.

Chris Oakley, ZOO Digital, UK Chris Oakley is the Chief Technology Officer (CTO) of ZOO Digital, a leading provider of cloud-based localization and digital distribution services for the global entertainment industry. With a career spanning over two decades in the technology and digital media sectors, Chris brings a wealth of experience and a visionary approach to his role at ZOO Digital. As CTO, Chris Oakley is responsible for overseeing the development and implementation of cutting-edge AI and ML technologies that power ZOO Digital's innovative services. Under his leadership, the company has continued to pioneer advancements in AI and ML cloud-based solutions, enabling efficient and scalable workflows for the localization and distribution of movies, TV shows, and other digital content.

Helena Moniz, President of EAMT & IAMT. University of Lisbon, Portugal. INESC-ID, Portugal

Helena Moniz is the President of the European Association for Machine Translation (2021-) and President of the International Association for Machine Translation (2023-). She is also the Vice-Coordinator of the Human Language Technologies Lab at INESC-ID, Lisbon. Helena is an Assistant Professor at the School of Arts and Humanities at the University of Lisbon, where she teaches Computational Linguistics, Computer Assisted Translation, and Machine Translation Systems and Post-editing. She is now in a very exciting project, coordinated by Unbabel, the Center for Responsible AI (<https://centerforresponsible.ai>), within the Portuguese Recovery and Resilience Plan, as Chair of the Ethics Committee. Helena graduated in Modern Languages and Literature at the School of Arts and Humanities, University of Lisbon (FLUL), in 1998. She took a Teacher Training graduation course in 2000, a Master's degree in Linguistics in 2007, and a PhD in Linguistics at FLUL in cooperation with the Technical University of Lisbon (IST) in 2013. She has been working at INESC-ID/CLUL since 2000, in several national and international projects involving multidisciplinary teams of linguists and speech processing engineers. Within these fruitful collaborations, she participated in more than 20 national and international projects. From 2015/09 to 2024/04, she was the PI of a bilateral project between INESC-ID and Unbabel, a translation company combining AI + post-editing, working on scalable Linguistic Quality Assurance processes for crowdsourcing. She was responsible for the implementation in 2015 of the MQM metric, the creation of the Linguistic Quality Assurance processes developed at Unbabel for Linguistic Annotation and Editors' Evaluation. She also worked on research projects, involving Linguistics, Translation, and Responsible AI, and products developed by the Labs Team, mostly cultural transcreation, high risk products, and silently controlled language metrics for dialogues. In a sentence, she is passionate about Language Technologies in a human-centric perspective and always feels like a child eager to learn!

Mirko Lorenz, Deutsche Welle, Germany Mirko Lorenz is an Innovation Manager working for Deutsche Welle, Germany's international broadcaster. He has been a member of the Research and Cooperation Team (ReCo) since 2008. One main outcome of his work is plain X, a 4-in-1 software to simplify content adaptation. In plain X, users can transcribe, translate, subtitle, and create (synthetic) voice-overs. Mirko has a master's in economics and history from the University of Cologne and a professional background in journalism. He co-founded Datawrapper, a tool to create charts and maps which is used in many large newsrooms worldwide.

Valter Mavrič, DG TRAD, European Parliament Valter Mavrič is Director-General of the Translation Service (DG TRAD) at the European Parliament (since 2016), where he was previously acting Director-General (from 2014), Director (from 2010) and Head of the Slovenian Translation Unit (from 2004). With an MA in applied linguistics and further training in translation, interpretation, linguistics and management, he has a long experience as manager, translator, interpreter and teacher of languages. He works in Slovenian, Italian, English, French, and Croatian and is currently preparing a PhD in strategic communication.

Moderator: Edward Gow-Smith, University of Sheffield, UK

Moderator: Carolina Scarton, University of Sheffield, UK

Table of Contents

Thesis Award	1
<i>Direct Speech Translation Toward High-Quality, Inclusive, and Augmented Systems</i> Marco Gaido	2
<i>Streaming Neural Speech Translation</i> Javier Iranzo-Sánchez	4
<i>Thesis: Model-based Evaluation of Multilinguality</i> Jannis Vamvas	6
Research: Technical	8
<i>Promoting Target Data in Context-aware Neural Machine Translation</i> Harritxu Gete and Thierry Etchegoyhen	9
<i>A Human Perspective on GPT-4 Translations: Analysing Faroese to English News and Blog Text Translations</i> Annika Simonsen and Hafsteinn Einarsson	24
<i>ReSeTOX: Re-learning attention weights for toxicity mitigation in machine translation</i> Javier García Gilabert, Carlos Escolano and Marta R. Costa-jussà	37
<i>Using Machine Translation to Augment Multilingual Classification</i> Adam King	59
<i>Recovery Should Never Deviate from Ground Truth: Mitigating Exposure Bias in Neural Machine Translation</i> Jianfei He, Shichao Sun, Xiaohua Jia and Wenjie Li	68
<i>Chasing COMET: Leveraging Minimum Bayes Risk Decoding for Self-Improving Machine Translation</i> Kamil Guttman, Mikołaj Pokrywka, Adrian Charkiewicz and Artur Nowakowski	80
<i>Mitra: Improving Terminologically Constrained Translation Quality with Backtranslations and Flag Diacritics</i> Iikka Hauhio and Théo Friberg	100
<i>Bootstrapping Pre-trained Word Embedding Models for Sign Language Gloss Translation</i> Euan McGill, Luis Chiruzzo and Horacio Saggion	116
<i>Quality Estimation with k-nearest Neighbors and Automatic Evaluation for Model-specific Quality Estimation</i> Tu Anh Dinh, Tobias Palzer and Jan Niehues	133
<i>SubMerge: Merging Equivalent Subword Tokenizations for Subword Regularized Models in Neural Machine Translation</i> Haiyue Song, Francois Meyer, Raj Dabre, Hideki Tanaka, Chenhui Chu and Sadao Kurohashi	147
<i>FAME-MT Dataset: Formality Awareness Made Easy for Machine Translation Purposes</i> Dawid Wisniewski, Zofia Rostek and Artur Nowakowski	164
<i>Iterative Translation Refinement with Large Language Models</i> Pinzhen Chen, Zhicheng Guo, Barry Haddow and Kenneth Heafield	181

<i>Detector–Corrector: Edit-Based Automatic Post Editing for Human Post Editing</i> Hiroyuki Deguchi, Masaaki Nagata and Taro Watanabe	191
<i>Assessing Translation Capabilities of Large Language Models involving English and Indian Languages</i> Vandan Mujadia, Ashok Urlana, Yash Bhaskar, Penumalla Aditya Pavani, Kukkapalli Shravya, Parameswari Krishnamurthy and Dipti Sharma	207
<i>Improving NMT from a Low-Resource Source Language: A Use Case from Catalan to Chinese via Spanish</i> Yongjian Chen, Antonio Toral, Zhijian Li and Mireia Farrús	229
<i>A Case Study on Context-Aware Neural Machine Translation with Multi-Task Learning</i> Ramakrishna Appicharla, Baban Gain, Santanu Pal, Asif Ekbal and Pushpak Bhattacharyya .	246
<i>Aligning Neural Machine Translation Models: Human Feedback in Training and Inference</i> Miguel Moura Ramos, Patrick Fernandes, António Farinhas and Andre Martins	258
<i>Enhancing Scientific Discourse: Machine Translation for the Scientific Domain</i> Dimitris Roussis, Sokratis Sofianopoulos and Stelios Piperidis	275
<i>Towards Tailored Recovery of Lexical Diversity in Literary Machine Translation</i> Esther Ploeger, Huiyuan Lai, Rik Van Noord and Antonio Toral	286
<i>Enhancing Gender-Inclusive Machine Translation with Neomorphemes and Large Language Models</i> Andrea Piergentili, Beatrice Savoldi, Matteo Negri and Luisa Bentivogli	300
Research: Translators & Users	315
<i>Prompting ChatGPT for Translation: A Comparative Analysis of Translation Brief and Persona Prompts</i> Sui He	316
<i>Exploring the Correlation between Human and Machine Evaluation of Simultaneous Speech Translation</i> Claudio Fantinuoli and Xiaoman Wang	327
<i>MTUncertainty: Assessing the Need for Post-editing of Machine Translation Outputs by Fine-tuning OpenAI LLMs</i> Serge Gladkoff, Lifeng Han, Gleb Erofeev, Irina Sorokina and Goran Nenadic	337
<i>Translators’ perspectives on machine translation uses and impacts in the Swiss Confederation: Navigating technological change in an institutional setting</i> Paolo Canavese and Patrick Cadwell	347
<i>Added Toxicity Mitigation at Inference Time for Multimodal and Massively Multilingual Translation</i> Marta R. Costa-jussà, David Dale, Maha Elbayad and Bokai YU	360
<i>LLMs in Post-Translation Workflows: Comparing Performance in Post-Editing and Error Analysis</i> Celia Soler Uguet, Fred Bane, Mahmoud Aymo, João Torres, Anna Zaretskaya and Tània Blanch Miró Blanch Miró	373
<i>Post-editors as Gatekeepers of Lexical and Syntactic Diversity: Comparative Analysis of Human Translation and Post-editing in Professional Settings</i> Lise Volkart and Pierrette Bouillon	387
<i>Exploring NMT Explainability for Translators Using NMT Visualising Tools</i> Gabriela Gonzalez-Saez, Mariam Nakhle, James Robert Turner, Fabien Lopez, Nicolas Ballier, Marco Dinarelli, Emmanuelle Esperança-Rodier, Sui He, Raheel Qader, Caroline Rossi, Didier Schwab and Jun Yang	396

<i>Mitigating Translationese with GPT-4: Strategies and Performance</i>	
Maria Kunilovskaya, Koel Dutta Chowdhury, Heike Przybyl, Cristina España-Bonet and Josef Van Genabith	411
<i>Translate your Own: a Post-Editing Experiment in the NLP domain</i>	
Rachel Bawden, Ziqian Peng, Maud Bénard, Éric Villemonte De La Clergerie, Raphaël Esamotunu, Mathilde Huguin, Natalie Kübler, Alexandra Mestivier, Mona Michelot, Laurent Romary, Lichao Zhu and François Yvon	431
<i>Pre-task perceptions of MT influence quality and productivity: the importance of better translator-computer interactions and implications for training</i>	
Vicent Briva-Iglesias and Sharon O’Brien	444
<i>Bayesian Hierarchical Modelling for Analysing the Effect of Speech Synthesis on Post-Editing Machine Translation</i>	
Miguel Rios, Justus Brockmann, Claudia Wiesinger, Raluca Chereji, Alina Secară and Dragoş Ciobanu	455
<i>Evaluation of intralingual machine translation for health communication</i>	
Silvana Deilen, Ekaterina Lapshinova-Koltunski, Sergio Hernandez Garrido, Julian Hörner, Christiane Maaß, Vanessa Theel and Sophie Ziemer	469
<i>Using Machine Learning to Validate a Novel Taxonomy of Phenomenal Translation States</i>	
Michael Carl, Sheng Lu and Ali Al-Ramadan	480
<i>Perceptions of Educators on MTQA Curriculum and Instruction</i>	
João Lucas Cavalheiro Camargo, Sheila Castilho and Joss Moorkens	492
<i>Comparative Quality Assessment of Human and Machine Translation with Best-Worst Scaling</i>	
Bettina Hiebl and Dagmar Gromann	507
<i>Quantifying the Contribution of MWEs and Polysemy in Translation Errors for English–Igbo MT</i>	
Adaeze Ngozi Ohuoba, Serge Sharoff and Callum Walker	537
<i>Analysis of the Annotations from a Crowd MT Evaluation Initiative: Case Study for the Spanish-Basque Pair</i>	
Nora Aranberri	548
Implementations & Case Studies	560
<i>A Case Study on Contextual Machine Translation in a Professional Scenario of Subtitling</i>	
Sebastian Vincent, Charlotte Prescott, Chris Bayliss, Chris Oakley and Carolina Scarton	561
<i>Training an NMT system for legal texts of a low-resource language variety South Tyrolean German - Italian</i>	
Antoni Oliver, Sergi Alvarez-Vidal, Egon Stemle and Elena Chiocchetti	573
<i>Implementing Gender-Inclusivity in MT Output using Automatic Post-Editing with LLMs</i>	
Mara Nunziatini and Sara Diego	580
<i>CantonMT: Cantonese to English NMT Platform with Fine-Tuned Models using Real and Synthetic Back-Translation Data</i>	
Kung Yin Hong, Lifeng Han, Riza Batista-Navarro and Goran Nenadic	590
<i>Advancing Digital Language Equality in Europe: A Market Study and Open-Source Solutions for Multilingual Websites</i>	
Andrejs Vasiljevs, Rinalds Vīksna, Neil Vacheva and Andis Lagzdīņš	600

<i>Exploring the Effectiveness of LLM Domain Adaptation for Business IT Machine Translation</i>	
Johannes Eschbach-Dymanus, Frank Essenberger, Bianka Buschbeck and Miriam Exel	610
<i>Creating and Evaluating a Multilingual Corpus of UN General Assembly Debates</i>	
Hannah Dorothy Bechara, Krishnamoorthy Manohara and Slava Jankin	623
<i>Generating subject-matter expertise assessment questions with GPT-4: a medical translation use-case</i>	
Diana Silveira, Marina Sánchez Torrón and Helena Silva Moniz	628
<i>Prompting Large Language Models with Human Error Markings for Self-Correcting Machine Translation</i>	
Nathaniel Berger, Stefan Riezler, Miriam Exel and Matthias Huck	636
<i>Estonian-Centric Machine Translation: Data, Models, and Challenges</i>	
Elizaveta Korotkova and Mark Fishel	647
Sponsors	661

Thesis Award

Direct Speech Translation Toward High-Quality, Inclusive, and Augmented Systems

Marco Gaido

Fondazione Bruno Kessler, Via Sommarive 18, 38123 Trento, Italy

University of Trento, Trento, Italy

Email: mgaido@fbk.eu, Phone: +39 3482670470

Supervisors: Marco Turchi, Matteo Negri

marco.turchi@zoom.us, negri@fbk.eu

When this PhD started, in November 2019, the translation of speech into text in a different language was mainly tackled with a cascade of automatic speech recognition (ASR) and machine translation (MT) models. However, a new paradigm was emerging, with the proposal of direct (or end-to-end) models designed to tackle the speech-to-text translation (ST) task in a single step. At that time, the main question within the ST community was: *will direct ST models be able to keep their promise and reach (or even outperform) the quality of cascade approaches?* Therefore, the initial phase of the PhD has been dedicated to building **high-quality** direct models, specifically under the practical scenario where lengthy audio files necessitate automated segmentation. The positive outcomes attained in terms of overall translation quality enabled the study of specific aspects of direct systems that are pivotal for meeting the real needs of end-users. Consequently, a significant portion of the PhD has been dedicated to analyzing and improving their behavior concerning two critical aspects: **inclusivity** (in terms of gender bias) and **augmented translation** (the integration of useful concepts and contextual information to help users' understanding). Below, I summarize the work I carried out on the above lines of research, and the related findings and achievements.

Translation Quality. Through the continuous experimentation of new techniques compared with the state of the art and evaluated in the challenging yearly international IWSLT evaluation campaign for speech translation, I contributed to closing the gap between the two paradigms, as attested by the first success of a direct system in the compe-

tion in 2020 (where the FBK model ranked 2nd, first among academic participants) and a thorough manual analysis carried out to compare the solutions (**ACL 2021**). Specifically, on one side I introduced training procedures and architectural solutions aimed at improving the translation quality of direct ST systems and their efficiency, reducing computational costs. On the other, I focused on how to limit the quality drops observed when the audio is not segmented according to a known reference but has to be automatically segmented into chunks processable by ST models.

As part of the first group of activities, I studied the best methods to transfer knowledge from an MT model into a direct ST system with knowledge distillation, highlighting not only the benefits but also its limitations, for which I provided an easy yet effective solution (**IWSLT 2020**). I also proposed a compression mechanism that leverages the prediction of a CTC module and dynamically reduces the length of the input sequence in the encoder of ST systems, improving both translation quality and computational efficiency (**EACL 2021**). Building on the CTC-compression module, I introduced Speechformer, the first architecture for direct ST that, enabled by an attention implementation with reduced computational complexity, avoids any fixed compression of the audio input, respecting the variability of the amount of information in speech signals and bringing significant quality gains (**EMNLP 2021**). Lastly, I showed the superfluity of the ASR pre-training when using an auxiliary CTC loss and the effectiveness of a simple data filtering procedure based on the transcript-to-translation character ratio (**IWSLT 2022**).

Moving to the second goal of coping with sub-optimal audio segmentation, I increased the robustness of direct ST models with regard to au-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

tomatic segmentation of the audio by fine-tuning them on resegmented training corpora and by providing the previous audio segment as contextual information (**Interspeech 2020**). Moreover, I proposed a new hybrid segmentation method that limits the quality degradation with respect to optimal segmentation based on the transcripts, which are unknown at inference time (**ICNLSP 2021**).

Inclusivity. Reckoning that a high overall quality is not enough to consider a technology ready for the users and driven by the ethical commitment and deep belief in the importance of raising awareness of the limitations – and even potential harms – of automatically-generated text in contemporary society, I devoted part of my PhD to studying the gender bias of direct ST systems. The goal was to ensure the fairness of automatic systems and equal opportunities for different groups of users to benefit from them. In this context, I disclosed how the pursuit of higher general performance can exacerbate gender representational disparities and proposed mitigation techniques that reduce the gender bias of ST models. To this aim, I explored different solutions to control the grammatical gender of words referred to the speaker (assuming that the gender of the speaker is known in advance), investigating for the first time the case in which the speakers’ gender conflicts with their vocal characteristics (**COLING 2020 Outstanding Paper**). In this context, I proposed automatic metrics tailored at disentangling the gender bias of a system from its overall quality, which has been validated through an extensive manual analysis, which also showed that ST models are nearly perfect in handling gender agreement and that the most biased part of speech is nouns (**ACL 2022**). Then, I unveiled the exacerbation of gender bias caused by a BPE segmentation of the target text in comparison with a character-based segmentation, and the proposal of a solution that goes beyond the trade-off between translation quality – BPE – and gender accuracy – char – (**ACL-Findings 2021**). Lastly, I demonstrated the increase in gender bias caused by distilling knowledge from MT and how to solve the issue with a simple fine-tuning (**CLiC-it 2020 Best Paper, IJCoL 2022**).

Augmented Translation. At last, motivated by the practical needs of interpreters and translators, my PhD evaluated the potential of direct ST systems in the “augmented translation” scenario, where the translation is enriched with contextual information

that eases its fruition. In particular, within the Smarter Interpreting¹ research project – aimed at the creation of to a new generation of computer-assisted interpreting (CAI) tools – the main focus was the translation and recognition of named entities (NEs), which constitute one of the most demanding challenges for interpreters. This strand of research activities started with the creation of a new benchmark (NEurRoparl-ST), used to assess the similar weaknesses of cascade and direct ST systems when it comes to NEs (**EMNLP 2021**). Having ascertained that person names are the most complex NE type for ST systems, I isolated the factors that contribute to this difficulty of ST systems (low frequency in the training data, names associated with languages not included in the source side of the training set) and proposed the adoption of multilingual models that jointly predict the transcript and the translation (giving more weight to the transcription) to mitigate such errors (**IWSLT 2023 Best Paper**). Moreover, in cases in which a dictionary of entities likely to appear in a given domain is available (a frequent condition in the interpreting sector), I showed that the accuracy of NEs (especially of person names) can be significantly improved by means of additional modules that first recognize which of them are present and then inject the corresponding translations as suggestions while generating the output (**ICASSP 2023**). The project was concluded by the introduction of models that jointly perform ST and NER, outperforming a pipeline of ST and NER systems while keeping the computational cost as low as that of a single direct ST model (**Interspeech 2023**).

Besides automatic evaluations on the proposed benchmark, the effectiveness of our solutions has been proved in two demos, carried out in April and December 2022, in which our joint ST and NER systems have been integrated into a new CAI tool that displays the translated NEs and domain-specific terminology in real time to the interpreter. In the first demo, students and professionals of the University of La Laguna performed a human-centric evaluation to assess the usefulness of the system for interpreters. The positive feedback of this analysis led to presenting the tool at international interpreting conferences,² where it has been introduced as the first 4th generation CAI system.

¹<https://smarter-interpreting.eu/> – financed by CDTI Neotec funds.

²https://ctn.hkbu.edu.hk/interpreting_conf2022/

Streaming Neural Machine Translation

Javier Iranzo-Sánchez
AppTek GmbH, València, Spain
jiranzo@apptek.com

Thesis Summary

Speech Translation (ST) is a subfield of Machine Learning (ML) that aims to automatically generate the text translation of a given audio waveform. Currently, the majority of the work in ST is concerned only with the offline task, that is, the task in which the entire input audio is available, and no real-time constraints exist. In contrast, in the online task the input audio is incrementally received as time passes, and the system must produce a translation of a partial input within a certain latency threshold, in a real-time fashion. Online ST is inherently a harder problem, because the partial input compromises the quality of the translation, and due to the need for real-time translation, the computational efficiency of the system cannot be ignored.

Traditionally, ST systems follow the cascade approach, in which the output of an Automatic Speech Recognition (ASR) system is fed into a Machine Translation (MT) system. Direct models are a more recent development, in which a single model receives the audio signal and generates the translation. The techniques presented on this thesis follow the cascade approach, but they can also be applied to the direct approach. Both approaches had achieved a similar level of performance at the time the thesis was written.

This thesis focuses on Streaming ST¹, a subtask of online ST in which the input is an unbounded audio stream. Streaming ST presents additional difficulties when compared with the standard online setup, and it is especially relevant because

many potential ST applications such as live lectures or simultaneous interpretation fall under the umbrella of streaming ST. The main goal of this thesis is to develop the tools and techniques that are required in order to create a working streaming ST solution. These are, specifically, a dataset for training and evaluating the ST models, a segmenter system that connects the output of the ASR system with the MT system, a streaming-ready evaluation metric and a streaming-specific MT model that can take advantage of contextual information.

The first challenge is the data scarcity problem faced when training ST systems. In order to alleviate this, a ST dataset is constructed starting from the official recordings of the proceedings of the European Parliament. The data is organized in triples, containing the audio jointly with its transcription and translation. It is a multilingual dataset with 10 different official European languages available both on the source and target side. Document-level information and metadata is included so that this dataset can be used for streaming ST.

The segmentation step is the next challenge to be addressed. The output of the streaming ASR system is a continuous stream of words, which needs to be segmented into semantically self-contained units to be translated by the MT system. We introduce a novel neural segmenter architecture, Direct Segmentation (DS), which considers the segmentation process as a classification problem. Using a sliding window approach, for every position of the ASR stream, the segmenter decides whether or not to produce a chunk by using a fixed local history and a small look-ahead window. The proposed architecture is computationally efficient while outperforming other segmentation approaches, and is able to work straight out

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹Streaming ST is also known in the literature as *long-form simultaneous ST*

of the box in the streaming scenario. Experiments are also performed showing that adding audio features to the segmenter improves performance. This work is then extended in order to evaluate the real latency for a simultaneous ST system that uses on-line ASR and MT systems as well as the proposed DS system. The results show how an acceptable translation quality can be reached at the same latency as a human interpreter (approximately 4 seconds).

The next challenge of streaming ST lies in how to actually evaluate the latency of the ST system under streaming conditions. This thesis introduces a novel evaluation procedure for streaming MT. Standard online MT metrics only work with short audio segments, evaluated in isolation, and do not take into account the sequential nature of the streaming scenario. Our proposed streaming evaluation method fixes these issues, and as a bonus, it can be applied to the standard metrics used for online MT with a small modification. Our proposal keeps track of a global latency score across the entire translation process, and uses a realignment step that matches translated words with the correct reference segment. A significant advantage of our proposal is that the evaluation procedure is not system/segmentation dependent and can be used to compare different systems, as well as maintaining the original interpretability of the metrics. Comparative experiments show that, unlike competing approaches, our proposal correctly ranks systems based on their latency, as well as keeping the previously mentioned properties.

Last but not least, we present a general methodology for building context-aware state-of-the-art streaming MT systems. This approach uses the insights developed in the previous publications in order to build a strong streaming baseline MT system, and improves it with a novel context-aware training methodology which obtains significant improvements. Further improvements are also obtained with a proposed Partial Bidirectional Encoder that has access to a larger portion of the input prefix. Our approach is similar to the concatenative approach used in context-aware MT, and uses a sliding window which contains the previous streaming history that has been produced during the translation process. History-augmented training samples are constructed from document-level corpora, and at inference time, the real streaming history is used. Extensive experiments show how

this approach achieves state-of-the-art results.

The full text of the thesis can be accessed at <https://doi.org/10.4995/Thesis/10251/199170>.

Supervisor Contact Details

Jorge Civera: Associate Professor, Universitat Politècnica de València, València, Spain. jorcisai@vrain.upv.es

Alfons Juan: Full Professor, Universitat Politècnica de València, València, Spain. ajuanci@vrain.upv.es

Acknowledgments

The author would like to publicly acknowledge the support received from his supervisors Jorge Civera and Alfons Juan before, during and after this PhD degree. Likewise, the author received many insightful comments from the thesis committee and external evaluators, consisting of Francisco Casacuberta Nolla, Jesús Andrés Ferrer, Marco Turchi, Felipe Sánchez Martínez and Marta Ruiz Costa-Jussà. This thesis has been developed with the financial support of the FPU scholarship program of the Government of Spain (FPU18/04135), EU's Horizon 2020 project X5gon (761758), as well as Government of Spain's Multisub (RTI2018-094879-B-I00) and Erasmus+ EXPERT (no. 20-226-093604-SCH) research projects. The author gratefully acknowledges the computer resources at Artemisa, funded by the European Union ERDF and Comunitat Valenciana as well as the technical support provided by the Instituto de Física Corpuscular, IFIC (CSIC-UV).

Model-based Evaluation of Multilinguality

Jannis Vamvas

Department of Computational Linguistics
University of Zurich
vamvas@cl.uzh.ch

The aim of this thesis was to extend the methodological toolbox for evaluating the ability of natural language processing systems to handle multiple languages. Neural machine translation (NMT) took the central role in this endeavor: NMT is inherently cross-lingual, and multilingual NMT systems, which translate from many source languages into many target languages, embody the concept of multilinguality in a very tangible way. In addition, NMT and specifically the perplexity of NMT systems can themselves be used as a tool for evaluating multilinguality.

Limitations of targeted evaluation methods for machine translation

In (Vamvas and Sennrich, 2021a), we identified a limitation of an existing targeted evaluation method, **contrastive evaluation using minimal pairs**. We discussed this limitation from a theoretical perspective by drawing a comparison between the conditions of contrastive evaluation and the concept of exposure bias.

We then performed experiments with English–German machine translation and demonstrated that testing implausible hypotheses using contrastive evaluation could lead to incorrect conclusions about the errors actually made by a system in practice. Finally, we proposed an effective mitigation approach, deriving minimal pairs from NMT-generated translations instead of human-written reference translations.

Contrastive conditioning: A novel approach to targeted evaluation

In (Vamvas and Sennrich, 2021b), we proposed **contrastive conditioning**, a novel targeted evaluation method for machine translation. Our idea is to analyze machine translations by measuring the perplexity of an “expert” NMT system that we provide with privileged information via a modified source sequence. Unlike some previous methods, contrastive conditioning can be used for a targeted evaluation of **black-box systems** such as commercial translation APIs. Another advantage of contrastive conditioning is that it requires few assumptions about the specific target language used, which allows for the scaling of automatic evaluation to many languages.

Two applications of contrastive conditioning

- In (Vamvas and Sennrich, 2021b), we used the method to quantify **overgeneralization bias** when translating ambiguous source expressions, which is a major challenge for machine translation. We hypothesized that lexical overgeneralization is more pronounced in NMT systems trained with **knowledge distillation**. Through the use of contrastive conditioning, we showed that distilled models are indeed more biased than non-distilled models, even if their overall quality is equal.
- In (Vamvas and Sennrich, 2022a), we demonstrated how contrastive conditioning can be applied to the automatic recognition of **erroneous omission and addition of content**. We performed a human evaluation study to validate our simple approach and found that the accuracy in detecting omission errors is comparable to that of a specialized quality estima-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

tion model that was trained on a large amount of synthetic data.

Translation cross-likelihood for semantic similarity

In the final publication included in the thesis (Vamvas and Sennrich, 2022b) we proposed a novel and robust way of using NMT perplexity for judging the similarity of sentence pairs, called **translation cross-likelihood**. We evaluated our approach on paraphrase identification and found that cross-likelihood tends to have a higher accuracy than previous approaches. We also found that translation-based similarity measures strongly outperform embedding-based measures in distinguishing between paraphrases and adversarial non-paraphrases. Finally, we highlighted the potential of evaluation based on NMT perplexity on the example of multilingual data-to-text generation.

Dissemination and Impact

A focus of this thesis has been the open sharing of research artifacts. All research code has been released on GitHub¹, including the NMTScore library² for computing translation perplexity. Whenever possible, open-source models and open datasets were used. Every paper was accompanied by a lay summary on the candidate’s research blog.³

Acknowledgments

The author would like to thank his Ph.D. supervisors, Rico Sennrich and Lena A. Jäger, and the doctoral committee members, Lena A. Jäger and Bill Byrne. The work presented in this thesis was funded by the Swiss National Science Foundation (project MUTAMUR; no. 176727).

Relevant Publications

Vamvas, Jannis and Rico Sennrich. 2021a. On the limits of minimal pairs in contrastive evaluation. In *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 58–68, Punta Cana, Dominican Republic, November. Association for Computational Linguistics.

Vamvas, Jannis and Rico Sennrich. 2021b. Contrastive conditioning for assessing disambiguation in MT: A

case study of distilled bias. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10246–10265, Online and Punta Cana, Dominican Republic, November. Association for Computational Linguistics.

Vamvas, Jannis and Rico Sennrich. 2022a. As little as possible, as much as necessary: Detecting over- and undertranslations with contrastive conditioning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 490–500, Dublin, Ireland, May. Association for Computational Linguistics.

Vamvas, Jannis and Rico Sennrich. 2022b. NMTScore: A multilingual analysis of translation-based text similarity measures. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 198–213, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.

¹<https://github.com/ZurichNLP>

²<https://github.com/ZurichNLP/nmtscore>

³<https://vamvas.ch>

Research: Technical

Promoting Target Data in Context-aware Neural Machine Translation

Harritsu Gete^{1,2*}

Thierry Etchegoyhen^{1*}

¹Vicomtech Foundation, Basque Research and Technology Alliance (BRTA)

²University of the Basque Country UPV/EHU

{hgete, tetchegoyhen}@vicomtech.org

Abstract

Standard context-aware neural machine translation (NMT) typically relies on parallel document-level data, exploiting both source and target contexts. Concatenation-based approaches in particular, still a strong baseline for document-level NMT, prepend source and/or target context sentences to the sentences to be translated, with model variants that exploit equal amounts of source and target data on each side achieving state-of-the-art results. In this work, we investigate whether target data should be further promoted within standard concatenation-based approaches, as most document-level phenomena rely on information that is present on the target language side. We evaluate novel concatenation-based variants where the target context is prepended to the source language, either in isolation or in combination with the source context. Experimental results in English-Russian and Basque-Spanish show that including target context in the source leads to large improvements on target language phenomena. On source-dependent phenomena, using only target language context in the source achieves parity with state-of-the-art concatenation approaches, or slightly underperforms, whereas combining source and target context on the source side leads to significant gains across the board.

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

*These authors contributed equally to this work.

1 Introduction

Significant progress has been achieved in Machine Translation within the Neural Machine Translation (NMT) paradigm (Sutskever et al., 2014; Bahdanau et al., 2015; Vaswani et al., 2017). For the most part though, most NMT models translate sentences in isolation, preventing the adequate translation on document-level phenomena such as cohesion, discourse coherence or intersentential anaphora resolution (Bawden et al., 2018; Läubli et al., 2018; Voita et al., 2019b; Lopes et al., 2020; Post and Junczys-Dowmunt, 2023). Among the various approaches to context-aware NMT, simple concatenation of context sentences, as initially proposed by Tiedemann and Scherrer (2017), remains a solid baseline typically used in practice with varying amounts of source-target context pairs (Agrawal et al., 2018; Junczys-Dowmunt, 2019; Majumder et al., 2022; Sun et al., 2022; Post and Junczys-Dowmunt, 2023).

Context-aware models typically rely on parallel document-level data, a scarce resource overall despite recent efforts to provide this type of resource (Barrault et al., 2019; Voita et al., 2019b; Gete et al., 2022). To the exception of approaches such as the monolingual repair framework of Voita et al. (2019a), context data in the source language is generally used as the core information to model context-awareness. However, most discourse-level phenomena feature information that is either present mainly in the target language (e.g., lexical cohesion, deixis) or in both the source and target languages (e.g., gender selection, ellipsis). Considering this, in this work we aim to explore the impact of promoting target language data in standard context-aware NMT.

Along these lines, we explore a simple

concatenation-based approach which consists in simply prepending context sentences from the target language to the source sentence to be translated, in isolation or in combination with source context. The underlying intuition is that contextual phenomena would be mainly modelled at the decoder level via target-side context information, whereas, on the encoder side, context data will be either ignored and copied, as foreign data, or also associated with source information to further model context. Using target language context data on the source side also enables the use of a standard NMT architecture and concatenation-based approach to context-aware NMT.

We show that replacing source context sentences with the target context already leads to significant gains for discourse-level phenomena that depend on target-language information, while achieving either parity or moderate degradation in contrastive accuracy on other phenomena. Combining both source and target context sentences on the source side leads to consistent significant improvements across the board. We establish our results on two language pairs, English-Russian and Basque-Spanish, for which contrastive test sets are publicly available on a range of phenomena that depend on the source and/or target language context.

In addition to accuracy results on specific phenomena, we compare the overall translation quality on parallel test sets as well. We also measure the impact of using either reference or machine-translated output as context at inference time, with only minor loss observed with the latter in our experiments. Finally, we evaluate the use of back-translated data, with similar comparative gains as those obtained using parallel document-level data. Overall, our experimental results indicate that promoting target context data within a standard NMT architecture can be a promising alternative for context-aware machine translation.

2 Related Work

One of the first methods proposed for document-level NMT is the concatenation of context sentences to the sentence to be translated, in either the source language only, or in both source and target languages (Tiedemann and Scherrer, 2017; Agrawal et al., 2018). This method does not require any architectural change and uses a fixed contextual window of sentences. It provides a robust baseline that often achieves performances

comparable to that of more sophisticated methods, in particular in high-resource scenarios (Lopes et al., 2020; Sun et al., 2022; Post and Junczys-Dowmunt, 2023). Variants of this approach include discounting the loss generated by the context (Lupo et al., 2022), extending model capacity (Majumder et al., 2022; Post and Junczys-Dowmunt, 2023) or encoding the specific position of the context sentences (Lupo et al., 2023; Gete and Etchegoyhen, 2023).

Alternative approaches include refining context-agnostic translations (Voita et al., 2019a; Mansimov et al., 2021) and modelling context information with specific NMT architectures (Jean et al., 2017; Li et al., 2020; Bao et al., 2021). More recently, the use of pretrained language models has been explored for the task, using them to encode the context (Wu et al., 2022), to initialize NMT models (Huang et al., 2023) or fusing the language model with a sentence-level translation model (Petrick et al., 2023). Directly using pretrained language models to perform translation can achieve competitive results, although these models might still produce critical errors and sometimes perform worse than conventional NMT models (Wang et al., 2023; Karpinska and Iyyer, 2023; Hendy et al., 2023).

Concatenation-based approaches vary regarding their use of context, exploiting either the source context (Zhang et al., 2018; Voita et al., 2018), the target context (Voita et al., 2019a) or both (Bawden et al., 2018; Agrawal et al., 2018; Xu et al., 2021; Majumder et al., 2022). The benefits of using context sentences in both the source and the target languages are also discussed in Müller et al. (2018), for a multi-encoder approach. Fernandes et al. (2021) conclude that concatenation-based models make more use of the target context than the source context, but Jin et al. (2023) show that the effectiveness of the target context versus the source context is highly dependent on the language pair involved. Close to the target-based approach we explore in this work, Scherrer et al. (2019) and Gete et al. (2023) include variants where target data is concatenated to the source sentence, notably showing that the target context is equally as important than source context, and particularly beneficial to address target-level phenomena. However, their experiments were limited to one target sentence, i.e. without prepending context on the target side. We show in this work that including the target con-

(a) Lexical cohesion: name translation
EN: Not for Julia. Julia has a taste for taunting her victims. RU: Не для Джулии[Julia]. Юлия*[Julia] умеет дразнить своих жертв.
(b) Deixis: register coherence
EU: Ez dago martetarrik zuen artean. Guztiak ari zarete ereduak lotu eta... ES: Ninguno de ustedes [form] es marciano. Todos vosotros estáis *[inf] siguiendo un modelo y... (None of you are Martians. You are all following a model and...)
(c) Gender selection
EU: Hori nire arreba da. Berak [?] zaindu zituen nire argazkiak. (That's my sister . He/She took care of my photos.) ES: Esa es mi hermana . Él * cuido mis fotos. (That's my sister . He * took care of my photos.)
(d) Verb phrase ellipsis
EN: Veronica, thank you, but you saw what happened. We all did [?]. RU: Вероника, спасибо, но ты видела, что произошло. Мы все хотели*. (Veronica, thank you, but you saw what happened. We all wanted * it.)

Table 1: Examples of document-level inconsistencies extracted from (Voita et al., 2019b) and (Gete et al., 2022).

text in both source and target languages is critical to achieve significant improvements overall.

Since standard NMT evaluation metrics such as BLEU (Papineni et al., 2002) are not well equipped to assess accuracy on discourse phenomena, several challenge test sets have been developed specifically to measure translations in context, via contrastive evaluations (Bawden et al., 2018; Müller et al., 2018; Voita et al., 2019b; Lopes et al., 2020; Nagata and Morishita, 2020; Gete et al., 2022; Currey et al., 2022). We include contrastive test sets that cover target-language phenomena such as deixis or lexical cohesion, as well as phenomena where the relevant context information is available in both source and target languages.

3 Exploiting Target Language Context

The main incentive for the promotion of target context data is the nature of the contextual phenomena of interest for machine translation, as these can be grouped into four broad categories depending on the location of the relevant contextual information.

In a first category would be discourse-level phenomena that require context information on the target language side, typically related to discursive cohesion in a broad sense (see examples *a* and *b* in Table 1). For instance, to maintain lexical cohesion beyond the sentence level, a quality translation should feature lexical repetition when necessary,

as it can mark emphasis or support question clarification. Another case is that of names with several possible translations, where translations must remain consistent throughout. Degrees of politeness and linguistic register in general also involve translation alternatives that are equally correct in isolation, but require consistency at the document level. In the case of pronouns, when the source antecedent has translation options in different grammatical genders, translation choices should be coherent throughout in the target language. In all of these cases, the relevant information involves previous translations in the target language.

In a second major category are phenomena for which the relevant context information is in both the source and the target context (examples *c* and *d* in Table 1). This includes word sense disambiguation scenarios, where different types of source or target elements may be relevant to perform disambiguation. Gender selection would also fall into this category, in those cases where translation options for the relevant contextual antecedent are unique or share the same gender. The resolution of elliptical constructions in the source language, with no equivalent in the target language, may also require context information from the source or the target language. Another instance for this type of phenomena would be the translation of Japanese zero pronouns into English (Nagata and Morishita,

(a)	ES: Hablé con mi amiga [fem]. Dijo que sí. EN: I talked to my friend [?]. She/He* said yes.
(b)	EN: You can't leave me! Don't go away! ES: ¡No puede dejarme! ¡No se vaya/te vayas* !

Table 2: Example of ambiguity where source context is necessary for disambiguation, in isolation (a) or in combination with the target context (b).

2020), where information on both sides can be relevant to determine the grammatical features of the target pronoun. Note that, even when contextual information is present in both the source and target languages, using source information for disambiguation can result in a lack of consistency in the target language, whenever incorrect translations are involved.¹

A third class of context-dependent phenomena exists, where source data are the only source of disambiguating information. This involves cases where the context includes the translation of a word marked for a specific category (e.g., gender) into an unmarked one, while the source sentence to be translated involves insufficient source information (e.g., a dropped pronoun) that needs to be translated into a marked element (e.g., a pronoun marked for gender). A typical example is provided in Table 2 *a*. In such a case, there would be insufficient information in the target language, as the proper translation of the dropped subject pronoun into *she* could only be determined from the gender of the source context antecedent *amiga* (*friend*).

Finally, a fourth broad category contains constructions where the source and target context need to be processed in combination for a correct translation. In the example *b* in Table 2, the source context subject *you* does not provide information about register, and neither does the target context in Spanish, since the verb *puede* can indicate either third person in informal register or second person in polite register. However, the source context indicates second person. Therefore combining both sources of context information, it can be derived that the translation should be second person in polite form.

Any target-only approach, such as monolingual repair (Voita et al., 2019a) or the target-only variant we also explore in this work, would only generate the correct translation in the latter

two classes of cases by either chance or training bias. Although these cases exist, it is unclear how widespread they actually are, compared to the other two main classes of contextual phenomena described above. In what follows, we set to compare the relative importance of source and target data across the main phenomena as represented in the selected document-level test suites.

4 Promoting Target Language Data

To explore the promotion of target language data, we simply prepend the target context sentences to the source sentence to be translated, either discarding or maintaining the source context sentences. On the target side, we evaluate the use of empty context as well as maintaining the target context sentences. We add a special token to separate the concatenated context sentences in all cases.

At inference time, in practice the previously translated sentences would be prepended as context. Since context translations can feature various degrees of correctness, we assess the approach under both ideal and average conditions. On parallel test sets, we measure the use of both correct reference context sentences (Section 6.1) and machine-translated ones (Section 8). On the contrastive test sets, only reference translations are used, as is standard practice, since target context coherence requirements prevent the use of non-reference context translations for fair evaluations (see the discussion in Section 8).

The prepended target-language data will need to be processed by the source language encoder under this approach, which might generate unwarranted noise. We hypothesise however that the encoder will essentially treat foreign language subwords as tokens to be copied directly into the target language, a typically simple operation for standard NMT models. We use BPE models jointly learned on merged source and target language data to facilitate this part of the process. Overall, the proposed approach provides the means to exploit target language data on the decoder side, without

¹Bawden et al. (2018) provide a contrastive test for these cases, where part of the source has been translated incorrectly but the translation is still required to be consistent overall.

any change to model architecture, while introducing data that might be easily processed via copying on the source side.

5 Experimental Setup

5.1 Data

We describe in turn below the datasets used to train and test our models. All selected datasets were normalised, tokenised and truecased using Moses (Koehn et al., 2007) and segmented with BPE (Sennrich et al., 2016), training a joint model over 32,000 operations. Tables 3 and 4 show corpora statistics for parallel and contrastive datasets respectively.

	EU-ES	EN-RU
TRAIN	1,753,726	6,000,000
DEV	3,051	10,000
TEST	6,078	10,000

Table 3: Parallel corpora statistics (number of sentences)

For Basque–Spanish, we selected the TANDO corpus (Gete et al., 2022), which contains parallel data from subtitles, news and literary documents. It includes two contrastive datasets for Basque to Spanish translation. The first one, GDR-SRC+TGT, centres on gender selection, with the disambiguating information present in both the source and target languages. The second one, COH-TGT, is meant to evaluate cases where, despite the absence in the source language of the necessary information to make a correct selection of gender or register, the translation must be contextually coherent using target-side information.

For English–Russian, we used the dataset described in Voita et al. (2019b), based on Open Sub-

EU-ES	Size	src	tgt	Dist.
GDR-SRC+TGT	300	✓	✓	≤ 5
COH-TGT	300		✓	≤ 5
EN-RU	Size	src	tgt	Dist.
Ellipsis infl.	500	✓	✓	≤ 3
Ellipsis VP	500	✓	✓	≤ 3
Deixis	2,500		✓	≤ 3
Lex. cohesion	1,500		✓	≤ 3

Table 4: Contrastive test sets: size (number of instances), required context information and distance to the disambiguating information (number of sentences)

titles excerpts (Lison et al., 2018). It includes 4 large-scale contrastive test sets for English to Russian translation. Two of these tests are related to ellipsis and contain the disambiguating information in both the source and target-side context: *Ellipsis infl.* assesses the selection of correct morphological noun phrase forms in cases where the source verb is elided, whereas *Ellipsis VP* evaluates the ability to predict the verb in Russian from an English sentence in which the verb phrase is elided. In the other two tests, the disambiguating information is only present in the target-side context: *Deixis* addresses politeness consistency in the target language, without nominal markers, whereas *Lexical Cohesion* focuses on the consistent translation of named entities in Russian.

5.2 Models

All models in our experiments are trained with Marian (Junczys-Downmunt et al., 2018) and rely on the Transformer-base architecture with the parameters described in Vaswani et al. (2017).

As a general baseline, we trained a sentence-level model using all source-target sentence pairs in the selected training datasets for each language pair. We then trained different variants of concatenation-based context-aware models, varying the type of context sentences prepended to the source and/or the target sentence, and adding a special token to separate the context.

We use the following convention to denote the models: *nton* uses the same amount of source and target data on each side, and represents the state-of-the-art baseline; *tgt-nton* uses target language data on both sides, discarding source context altogether; *nto1* and *tgt-nto1* are variants of the previous models that use no context sentences in the target language; finally, *src+tgt-nton* and *tgt+src-nton* are variants where target context sentences are combined with source context sentences, by prepending them after or before the latter, respectively. For convenience, we will refer to the *tgt-nton*, *src+tgt-nton* and *tgt+src-nton* variants as *X-tgt-nton*, as they share the use of target context on both sides. In Appendix A, we provide a diagram to illustrate data composition for each model.

Given the size of the context for each dataset, we have $n=6$ for Basque–Spanish models and $n=4$ for English–Russian models. All context-aware models were initialised with the weights of the sentence-level baseline.

Note that we discarded *1ton* models, as they present two main challenges. Within a standard concatenation approach, we would be tasking the model to learn a transformation from a single source sentence to both the context and the target sentence, although the target context cannot be derived from the source sentence, obviously. Alternatively, a *1ton* model could be designed via changes in the NMT architecture, with forced decoding over the specified target context at both training and inference time. The required architectural changes were beyond the scope of this work, although this type of model might be worth exploring in more details.

6 Results

6.1 Parallel Tests

We first compared models in terms of BLEU on the parallel test sets, using SacreBLEU (Post, 2018)². Statistical significance was computed via paired bootstrap resampling (Koehn, 2004), for $p < 0.05$.³ The results are shown in Table 5.

In Basque–Spanish, the *nton*, *tgt-nton*, and *src+tgt-nton* models performed better than the alternatives, with no statistically significant differences between the three, with the *tgt+src-nton* achieving slightly lower results. All three were notably significantly better than the baseline and the models which used only a single reference in the target language. In English–Russian, all *X-tgt-nton* model variants, that included target context data on the source side, outperformed all other models, including the standard *nton* model.

	EU-ES	EN-RU
Sentence-level	31.20	31.09
<i>nto1</i>	29.91	31.48
<i>tgt-nto1</i>	29.43	31.03
<i>nton</i>	31.96	31.20
<i>tgt-nton</i>	31.82	32.29
<i>src+tgt-nton</i>	31.94	32.32
<i>tgt+src-nton</i>	31.56	32.49

Table 5: BLEU results on the parallel test sets.

Sentence-level metrics are typically insufficient to assess translation quality at the document level (Wong and Kit, 2012), and conclusions should not

²nrefs:1lcase:mixedlff:nltk:13alsmooth:explversion:2.3.1

³In all tables, best scores given the statistical test at hand are shown in bold.

be drawn from the above results regarding context-aware ability of the different models. They do however indicate several tendencies at the sentence level. First, the proposed use of target context data on both sides was not detrimental in terms of translation quality, as the *X-tgt-nton* models performed on a par with, or better than, the other variants. Secondly, the lower results obtained by the *nto1* and *tgt-nto1* models seem to indicate that (i) removing target context data on the decoder side can be detrimental, as in EU-ES, and (ii) using source or target language data on the encoder side can lead to similar BLEU results, as was the case in both language pairs.

Note that the results above were obtained with reference translations, in an ideal scenario where the context is correctly translated. In Section 8, we present additional results using machine-translated context, to measure the impact of eventual errors in target context translation.

6.2 Challenge Tests

We evaluated the different models on the challenge test sets both in terms of BLEU and in terms of accuracy of the contrastive evaluation. Statistical significance of accuracy results was computed using McNemar’s test (McNemar, 1947), for $p < 0.05$. The results are shown in Tables 6 and 7.

Considering both language pairs, the first notable results are the significant gains achieved by the *src+tgt-nton* and *tgt+src-nton* models, which outperformed all other variants overall, in terms of both BLEU scores and contrastive accuracy. The *tgt-nton* model, where source context was discarded altogether, also outperformed the baselines in terms of BLEU in all but one case, and either matched the other two target-based variants in half of the scenarios, or was outperformed by these variants in the other three cases. In terms of contrastive accuracy, it also outperformed the baselines by a wide margin on target-oriented phenomena while achieving parity or resulting in accuracy loss on other phenomena. Overall, the best performing and most consistent variant across datasets and metrics was the *src+tgt-nton* variant.

On all target-related phenomena, the *X-tgt-nton* models outperformed all alternatives, and in particular the standard *nton* variant by large margins. In terms of accuracy, in EU-ES on the COH-TGT test, the *tgt-nton* model already outperformed the baseline by 27.67 points and the *nton* model by

	GDR-SRC+TGT		COH-TGT	
	BLEU	ACC.	BLEU	ACC.
Sentence-level	36.28	53.67	35.04	54.00
<i>nto1</i>	36.82	66.33	33.23	53.00
<i>tgt-nto1</i>	36.79	66.33	37.31	74.00
<i>nton</i>	40.45	77.67	35.89	65.33
<i>tgt-nton</i>	39.05	72.67	39.61	81.67
<i>src+tgt-nton</i>	41.29	78.67	40.23	84.67
<i>tgt+src-nton</i>	42.35	78.67	39.86	82.67

Table 6: BLEU and accuracy results on the Basque–Spanish challenge tests.

	Ellipsis infl.		Ellipsis VP		Deixis		Lex. Cohesion	
	BLEU	ACC.	BLEU	ACC.	BLEU	ACC.	BLEU	ACC.
Sentence-level	30.81	51.80	22.20	27.80	28.10	50.04	31.52	45.87
<i>nto1</i>	32.69	54.60	30.24	65.40	28.20	50.04	29.47	45.87
<i>tgt-nto1</i>	32.28	53.60	23.59	29.00	28.30	50.56	30.37	45.87
<i>nton</i>	36.97	75.20	29.59	62.60	27.15	82.48	27.89	45.93
<i>tgt-nton</i>	40.69	70.00	30.75	60.00	34.17	87.48	30.98	49.47
<i>src+tgt-nton</i>	40.98	77.20	35.84	77.60	34.38	87.48	31.75	53.07
<i>tgt+src-nton</i>	42.02	75.60	34.46	74.88	34.07	88.28	31.33	51.00

Table 7: BLEU and accuracy results in English–Russian challenge tests.

16.34 points, with even higher accuracy gains for the best-performing *src+tgt-nton* model (+19.34). In EN-RU, on *Deixis* gains of up to 38.24 and 5.8 points were achieved against the baseline and *nton* model, respectively; on the *Lexical Cohesion* test set, the gains reached 7.2 and 7.14 points, respectively. On these target-oriented test-sets, all X-*tgt-nton* model also achieved comparable gains in terms of BLEU scores, with a maximum against the *nton* model of +4.34 points in EU-ES, +7.23 in EN-RU on *Deixis*, and +3.86 in EN-RU on the *Lexical cohesion* test.

Turning now to the test sets where relevant context information is available in either both the source and target languages, or perhaps only in the source language in some cases, the results are more balanced between the *nton* baseline and the X-*tgt-nton* variants, although the *src+tgt-nton* achieved the best results overall in terms of both BLEU and accuracy. On *Ellipsis VP*, the latter notably achieved gains of 15 accuracy points, with the *tgt+src-nton* variant a close second at +12.28. On *Ellipsis infl.* and GDR-SRC-TGT, the gains were more limited, with a maximum of +1 and +2 accuracy points for the *src+tgt-nton* model against the *nton* baseline, respectively, although signifi-

cant BLEU gains of up to +3.3 and +5.05 were observed on these test sets, respectively.

Unsurprisingly, on these three datasets where source information is a relevant factor, in combination or in isolation, the *tgt-nton* model underperformed, though in accuracy only and to a limited extent on *Ellipsis VP*, for instance. This variant also significantly outperformed the *nton* baseline in terms of BLEU on *Ellipsis infl.*, with a 3.60 points gain. To further determine the impact of source and target context and more precisely assess the limits of this type of model, more fine-grained challenge tests would be needed to distinguish between cases that can solely be resolved with source context information and those where either side of context provides sufficient information.

Regarding the other two contextual variants, *nto1* and *tgt-nto1*, which used no context information on the target side of the input, the results in accuracy were similar overall, performing on a par with the sentence-level baseline on *Lexical Cohesion*, *Deixis* and COH-TGT for *nto1*. This was expected for the *nto1* models, as the relevant information is in the target language in these cases, which these models have no access to.

Overall, promoting target data in a

concatenation-based approach achieved large improvements across the board over the sentence-level and *nton* baselines. Replacing source context data altogether with the target context already improved significantly on target-context phenomena, while achieving relatively close results in the other cases. Combining source and target context provided the best balance however, achieving the best results in all cases. In particular, the *src+tgt-nton* proved optimal and we discarded the slightly worse *tgt+src-nton* variant in the remainder of this work.

7 Using Back-translated Data

When document-level parallel data are lacking, monolingual data in the target language can be exploited within concatenation-based approaches via back-translation (Junczys-Dowmunt, 2019; Sugiyama and Yoshinaga, 2019; Huo et al., 2020). Some level of degradation is expected, depending on the quality of the model used to back-translate the target data, and we also expect the models to be impacted differently: the target sentence and its back-translation would be identical for all models, as would be the original target context sentences, but the *nton* and the *src+tgt-nton* models also require back-translated context, unlike the *tgt-nton* model.

For comparison purposes we back-translated the target side of the training data for both language pairs, using a sentence-level model trained on the parallel data, and trained the main model variants strictly on the back-translated data.⁴ The results are shown in Tables 8, 9 and 10, contrasting the use of parallel (PA) and back-translated (BT) data.

The overall degradation using BT data was more salient in EU-ES than in EN-RU, which is likely due to the differences in training data size and the resulting quality of the respective models. In both cases, the *X-tgt-nton* variants proved more robust than the *nton* model. This is likely due to the latter having as context only the back-translation of the target context, while the former contain, alone or in combination with the back-translation, the original target context.

Overall, the tendencies observed using parallel data were replicated with back-translated data, with the *src+tgt-nton* model being the top-

⁴Note that we did not mix back-translated data with the original parallel data, to strictly contrast the approaches in their ability to exploit monolingual back-translated data.

	EU-ES	EN-RU
Sentence-level (PA)	31.20	31.09
<i>nton</i> (PA)	31.96	31.20
<i>tgt-nton</i> (PA)	31.82	32.29
<i>src+tgt-nton</i> (PA)	31.94	32.32
<i>nton</i> (BT)	25.46	29.21
<i>tgt-nton</i> (BT)	27.33	30.10
<i>src+tgt-nton</i> (BT)	31.27	29.39

Table 8: BLEU results on the parallel test sets using parallel (PA) and back-translated (BT) data.

performing variant across the board, and the *tgt-nton* a close second on target-context phenomena but performing worse than the *nton* model in accuracy on the GDR-SRC+TGT and *Ellipsis infl.* with BT data. Perhaps more surprising are the results achieved by the *src+tgt-nton* model, trained on BT data, on the *Lexical cohesion* test set, where it outperformed the same variant trained on parallel data by 13 points. Additional datasets might be warranted to further assess the tendencies for these models, but the results on the available datasets in terms of accuracy seem to indicate that the use of BT data is viable, and particularly exploitable by the *X-tgt-nton* models overall. We conjecture that this is mainly due to the fact that these approaches promote target language data which are in essence correct, while discarding or reducing the role of source context data which are likely to feature back-translation errors.

8 Machine-translated Target Context

Following standard practice, so far we used the reference target context instead of the machine-translated output in our evaluations. This is meant to remove potential noise in terms of context translation errors and evaluate the approaches on their translation accuracy given a correct context. Using reference translations also allows for an evaluation of phenomena where more than one context translation would be correct – e.g. *box* translated as *boîte* (fem.) instead of *carton* (masc.) in French – but the contrastive evaluation relies on one of these translations being selected and contextual phenomena, such as coherence, are evaluated accordingly. A correct but different context translation would unfairly affect the evaluation.

Still, in practice, at inference time there are no reference translations, of course. Whereas X-to1

	GDR-SRC+TGT		COH-TGT	
	BLEU	ACC.	BLEU	ACC.
Sentence-level	36.28	53.67	35.04	54.00
<i>nton</i> (PA)	40.45	77.67	35.89	65.33
tgt- <i>nton</i> (PA)	39.05	72.67	39.61	81.67
src+tgt- <i>nton</i> (PA)	41.25	78.67	40.23	84.67
<i>nton</i> (BT)	41.58	76.00	31.02	67.00
tgt- <i>nton</i> (BT)	40.22	74.00	34.62	81.33
src+tgt- <i>nton</i> (BT)	45.67	77.33	42.67	84.67

Table 9: Results on Basque–Spanish contrastive tests with parallel (PA) and back-translated (BT) data.

	Ellipsis infl.		Ellipsis VP		Deixis		Lex. cohesion	
	BLEU	ACC.	BLEU	ACC.	BLEU	ACC.	BLEU	ACC.
Sentence-level	30.81	51.80	22.20	27.80	28.10	50.04	31.52	45.87
<i>nton</i> (PA)	36.97	75.20	29.59	62.60	27.15	82.48	27.89	45.93
tgt- <i>nton</i> (PA)	40.69	70.00	30.75	60.00	34.17	87.48	30.98	49.47
src+tgt- <i>nton</i> (PA)	40.98	77.20	35.84	77.60	34.38	87.48	31.75	53.07
<i>nton</i> (BT)	35.63	78.60	28.84	69.40	25.66	83.92	28.29	46.20
tgt- <i>nton</i> (BT)	39.25	73.60	31.86	57.60	31.84	87.84	29.81	49.20
src+tgt- <i>nton</i> (BT)	41.96	81.20	35.23	76.00	31.63	87.36	31.68	66.07

Table 10: Results on English–Russian contrastive tests with parallel (PA) and back-translated (BT) data.

	EU-ES	EN-RU
Sentence-level	31.20	31.09
<i>nton</i>	31.96	31.20
tgt- <i>nton</i> (RF)	31.82	32.29
tgt- <i>nton</i> (MT)	31.08	31.52
src+tgt- <i>nton</i> (RF)	31.94	32.32
src+tgt- <i>nton</i> (MT)	30.93	31.31

Table 11: BLEU results on the parallel test sets using reference (RF) and machine-translated (MT) context.

model should not be impacted at all, the X-tgt-*nton* models are susceptible to suffer from errors in the translation of the context. To measure this aspect, we computed BLEU scores using machine-translated target sentences for X-tgt-*nton* models. The results are shown in Table 11.

Using MT output resulted in a slight degradation for EU-ES, with results on a par with the sentence-level baseline and at most 1.01 points loss compared to the use of reference translations. For EN-RU, all models achieved comparable results except those that relied on reference translations, with

gains of approximately 1 point for the latter. As previously noted, the BLEU metric is known to be deficient for context-aware model evaluation, and contrastive tests provide more precise benchmarks. However, measuring MT context in terms of contrastive accuracy is not a valid option, as challenge tests rely on specific context translation choices, and the reference context is provided instead in standard practice. Note that *nton* models would also be impacted in terms of contrastive accuracy, since MT output would affect decoding.⁵

Evaluating approaches based on promoting target data in a practical scenario with imperfect machine-translated context thus faces important limitations with current document-level evaluation protocols. A proper assessment of the impact of machine-translated context would also need to take into account the quality of the translation model itself, with larger models expected to minimise context translation errors in this type of approach. We leave these aspects for future research.

⁵For completeness, in Appendix B we provide results in terms of BLEU and accuracy on the challenge tests using machine-translated context.

<i>Dist</i>	% cases	GDR-SRC+TGT			% cases	COH-TGT		
		<i>nton</i>	<i>tgt-nton</i>	<i>src+tgt-nton</i>		<i>nton</i>	<i>tgt-nton</i>	<i>src+tgt-nton</i>
1	64.67%	77.32	70.10	76.80	62.34%	69.52	85.03	86.10
2	20.67%	91.23	85.48	85.48	20.67%	66.13	90.32	85.48
3	9.33%	72.41	71.43	71.43	9.67%	51.72	72.41	75.86
4	2.00%	57.14	57.14	85.71	6.00%	50.00	83.33	83.33
5	3.33%	66.67	55.56	88.89	1.33%	25.00	50.00	75.00

Table 12: Accuracy results in Basque–Spanish according to relevant context distance.

<i>Dist</i>	% cases	Deixis			% cases	Lex. Cohesion		
		<i>nton</i>	<i>tgt-nton</i>	<i>src+tgt-nton</i>		<i>nton</i>	<i>tgt-nton</i>	<i>src+tgt-nton</i>
1	33.33%	88.66	90.49	89.63	42.75%	46.27	51.45	57.53
2	33.33%	85.82	90.07	91.02	31.50%	45.87	47.39	50.00
3	33.33%	73.02	81.89	81.77	25.75%	45.43	48.56	49.09

Table 13: Accuracy results in English–Russian according to relevant context distance.

9 Accuracy At Distance

The results so far were measured considering context as a whole. To achieve a more fine-grained view of the differences between approaches, we computed their accuracy in terms of the distance between the current sentence and the disambiguating context information, expressed in number of sentences. The results are shown in Tables 12 and 13, indicating the distance and the percentages of cases in the corresponding dataset.

The main observable tendency is that of the decreasing accuracy over distance for the *nton* model, in all cases but GDR-SRC+TGT at distance 2 (where all models perform better), in contrast with the significantly more robust accuracy of the *src+tgt-nton* model at larger distances, for Basque-Spanish in particular. The *tgt-nton* model exhibits mixed tendencies, improving or maintaining accuracy over distance 1 in some cases, but also degrading at larger distances (GDR-SRC+TGT or COH-TGT, at *dist*=5). Note though that larger distances are under-represented in the Basque-Spanish test sets, and may thus not be as representative.

10 Conclusions

In this work, we investigated the promotion of target context data within a standard concatenation-based approach to context-aware neural machine translation. The main incentive revolves around the fact that, for most contextual phenomena of interest for document-level machine translation, the relevant information is either in the target language

or distributed on the source and target sides.

We studied simple model variants where target context sentences are concatenated to the source sentence, either in isolation or in combination with the source context. Our results in Basque-Spanish and English-Russian, over five datasets showcasing different types of contextual phenomena, showed large improvements in terms of contrastive accuracy and BLEU scores. Models where the source context was discarded altogether achieved parity or slightly underperformed on phenomena involving both source and target contexts. The variants based on augmenting the source context with target data achieved the best results across the board and were also shown to be more accurate in handling context at larger distances.

We further evaluated the use of back-translated data, with models merging target and source matching or outperforming variants trained on parallel data. We also measured the impact of using machine-translated context, although only in a limited way given current evaluation protocols for context-aware models, with slight degradation observed in terms of BLEU. The use of more robust baseline models, trained on larger volumes of data, could mitigate the observed effects.

The proposed approach promoting target data requires no changes to the standard NMT architecture and provides significant gains over strong baselines. Although it also implies larger contexts when merging source and target context, it might be worth further exploring this type of approach and the respective roles of source and target context data in neural machine translation.

References

- Agrawal, Ruchit, Marco Turchi, and Matteo Negri. 2018. Contextual handling in neural machine translation: Look behind, ahead and on both sides. In *Proceedings of the 21st Annual Conference of the European Association for Machine Translation (EAMT)*.
- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*.
- Bao, Guangsheng, Yue Zhang, Zhiyang Teng, Boxing Chen, and Weihua Luo. 2021. G-transformer for document-level machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3442–3455, Online, August. Association for Computational Linguistics.
- Barrault, Loïc, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. Findings of the 2019 conference on machine translation (WMT19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Florence, Italy, August. Association for Computational Linguistics.
- Bawden, Rachel, Rico Sennrich, Alexandra Birch, and Barry Haddow. 2018. Evaluating discourse phenomena in neural machine translation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1304–1313, New Orleans, Louisiana, June. Association for Computational Linguistics.
- Currey, Anna, Maria Nadejde, Raghavendra Reddy Pappagari, Mia Mayer, Stanislas Lauly, Xing Niu, Benjamin Hsu, and Georgiana Dinu. 2022. MT-GenEval: A counterfactual and contextual dataset for evaluating gender accuracy in machine translation. In Goldberg, Yoav, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4287–4299, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.
- Fernandes, Patrick, Kayo Yin, Graham Neubig, and André F. T. Martins. 2021. Measuring and increasing context usage in context-aware machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6467–6478, Online, August. Association for Computational Linguistics.
- Gete, Harritxu and Thierry Etchegoyhen. 2023. An evaluation of source factors in concatenation-based context-aware neural machine translation. In Mitkov, Ruslan and Galia Angelova, editors, *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing*, pages 399–407, Varna, Bulgaria, September. IN-COMA Ltd., Shoumen, Bulgaria.
- Gete, Harritxu, Thierry Etchegoyhen, David Ponce, Gorka Labaka, Nora Aranberri, Ander Corral, Xabier Saralegi, Igor Ellakuria, and Maite Martin. 2022. TANDO: A corpus for document-level machine translation. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3026–3037, Marseille, France.
- Gete, Harritxu, Thierry Etchegoyhen, and Gorka Labaka. 2023. What works when in context-aware neural machine translation? In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 147–156, Tampere, Finland, June. European Association for Machine Translation.
- Hendy, Amr, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation.
- Huang, Zhihong, Longyue Wang, Siyou Liu, and Derek F. Wong. 2023. How does pretraining improve discourse-aware translation?
- Huo, Jingjing, Christian Herold, Yingbo Gao, Leonard Dahlmann, Shahram Khadivi, and Hermann Ney. 2020. Diving deep into context-aware neural machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 604–616, Online, November. Association for Computational Linguistics.
- Jean, Sebastien, Stanislas Lauly, Orhan Firat, and Kyunghyun Cho. 2017. Neural machine translation for cross-lingual pronoun prediction. In *Proceedings of the Third Workshop on Discourse in Machine Translation*, pages 54–57, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Jin, Linghao, Jacqueline He, Jonathan May, and Xuezhe Ma. 2023. Challenges in context-aware neural machine translation. In Bouamor, Houde, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15246–15263, Singapore, December. Association for Computational Linguistics.

- Junczys-Dowmunt, Marcin, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia, July. Association for Computational Linguistics.
- Junczys-Dowmunt, Marcin. 2019. Microsoft translator at WMT 2019: Towards large-scale document-level neural machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 225–233, Florence, Italy, August. Association for Computational Linguistics.
- Karpinska, Marzena and Mohit Iyyer. 2023. Large language models effectively leverage document-level context for literary translation, but critical errors persist.
- Koehn, Philipp, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, pages 177–180, Prague, Czech Republic.
- Koehn, Philipp. 2004. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain, July. Association for Computational Linguistics.
- Läubli, Samuel, Rico Sennrich, and Martin Volk. 2018. Has machine translation achieved human parity? a case for document-level evaluation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4791–4796, Brussels, Belgium, October–November. Association for Computational Linguistics.
- Li, Bei, Hui Liu, Ziyang Wang, Yufan Jiang, Tong Xiao, Jingbo Zhu, Tongran Liu, and Changliang Li. 2020. Does multi-encoder help? a case study on context-aware neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3512–3518, Online, July. Association for Computational Linguistics.
- Lison, Pierre, Jörg Tiedemann, and Milen Kouylekov. 2018. OpenSubtitles2018: Statistical rescored of sentence alignments in large, noisy parallel corpora. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May. European Language Resources Association (ELRA).
- Lopes, António, M. Amin Farajian, Rachel Bawden, Michael Zhang, and André F. T. Martins. 2020. Document-level neural MT: A systematic comparison. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 225–234, Lisboa, Portugal, November. European Association for Machine Translation.
- Lupo, Lorenzo, Marco Dinarelli, and Laurent Besacier. 2022. Focused concatenation for context-aware neural machine translation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 830–842, Abu Dhabi, United Arab Emirates (Hybrid), December.
- Lupo, Lorenzo, Marco Dinarelli, and Laurent Besacier. 2023. Encoding sentence position in context-aware neural machine translation with concatenation. In *The Fourth Workshop on Insights from Negative Results in NLP*, pages 33–44, Dubrovnik, Croatia, May.
- Majumder, Suvodeep, Stanislas Lauly, Maria Nadejde, Marcello Federico, and Georgiana Dinu. 2022. A baseline revisited: Pushing the limits of multi-segment models for context-aware translation.
- Mansimov, Elman, Gábor Melis, and Lei Yu. 2021. Capturing document context inside sentence-level neural machine translation models with self-training. In *Proceedings of the 2nd Workshop on Computational Approaches to Discourse*, pages 143–153, Punta Cana, Dominican Republic and Online, November. Association for Computational Linguistics.
- McNemar, Quinn. 1947. Note on the sampling error of the difference between correlated proportions or percentages. *Psychometrika*, 12:153–157.
- Müller, Mathias, Annette Rios, Elena Voita, and Rico Sennrich. 2018. A large-scale test set for the evaluation of context-aware pronoun translation in neural machine translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 61–72, Brussels, Belgium, October. Association for Computational Linguistics.
- Nagata, Masaaki and Makoto Morishita. 2020. A test set for discourse translation from Japanese to English. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 3704–3709, Marseille, France, May. European Language Resources Association.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Petrick, Frithjof, Christian Herold, Pavel Petrushkov, Shahram Khadivi, and Hermann Ney. 2023. Document-level language models for machine translation. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings*

- of the *Eighth Conference on Machine Translation*, pages 375–391, Singapore, December. Association for Computational Linguistics.
- Post, Matt and Marcin Junczys-Dowmunt. 2023. Escaping the sentence-level paradigm in machine translation. *arXiv preprint arXiv:2304.12959v1*.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October. Association for Computational Linguistics.
- Scherrer, Yves, Jörg Tiedemann, and Sharid Loáigiga. 2019. Analysing concatenation approaches to document-level NMT in two different domains. In *Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019)*, pages 51–61, Hong Kong, China, November. Association for Computational Linguistics.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August. Association for Computational Linguistics.
- Sugiyama, Amane and Naoki Yoshinaga. 2019. Data augmentation using back-translation for context-aware neural machine translation. In *Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019)*, pages 35–44, Hong Kong, China, November. Association for Computational Linguistics.
- Sun, Zewei, Mingxuan Wang, Hao Zhou, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Lei Li. 2022. Rethinking document-level neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3537–3548, Dublin, Ireland, May. Association for Computational Linguistics.
- Sutskever, Ilya, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112.
- Tiedemann, Jörg and Yves Scherrer. 2017. Neural machine translation with extended context. In *Proceedings of the Third Workshop on Discourse in Machine Translation*, pages 82–92, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems 30*, pages 5998–6008.
- Voita, Elena, Pavel Serdyukov, Rico Sennrich, and Ivan Titov. 2018. Context-aware neural machine translation learns anaphora resolution. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1264–1274, Melbourne, Australia, July. Association for Computational Linguistics.
- Voita, Elena, Rico Sennrich, and Ivan Titov. 2019a. Context-aware monolingual repair for neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 877–886, Hong Kong, China, November. Association for Computational Linguistics.
- Voita, Elena, Rico Sennrich, and Ivan Titov. 2019b. When a good translation is wrong in context: Context-aware machine translation improves on deixis, ellipsis, and lexical cohesion. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1198–1212, Florence, Italy, July. Association for Computational Linguistics.
- Wang, Longyue, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-level machine translation with large language models.
- Wong, Billy T. M. and Chunyu Kit. 2012. Extending machine translation evaluation metrics with lexical cohesion to document level. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 1060–1068, Jeju Island, Korea, July. Association for Computational Linguistics.
- Wu, Xueqing, Yingce Xia, Jinhua Zhu, Lijun Wu, Shufang Xie, and Tao Qin. 2022. A study of BERT for context-aware neural machine translation. *Mach. Learn.*, 111(3):917–935.
- Xu, Mingzhou, Liangyou Li, Derek F. Wong, Qun Liu, and Lidia S. Chao. 2021. Document graph for neural machine translation. In Moens, Marie-Francine, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 8435–8448. Association for Computational Linguistics.
- Zhang, Jiacheng, Huanbo Luan, Maosong Sun, Feifei Zhai, Jingfang Xu, Min Zhang, and Yang Liu. 2018. Improving the transformer translation model with document-level context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 533–542, Brussels, Belgium, October-November. Association for Computational Linguistics.

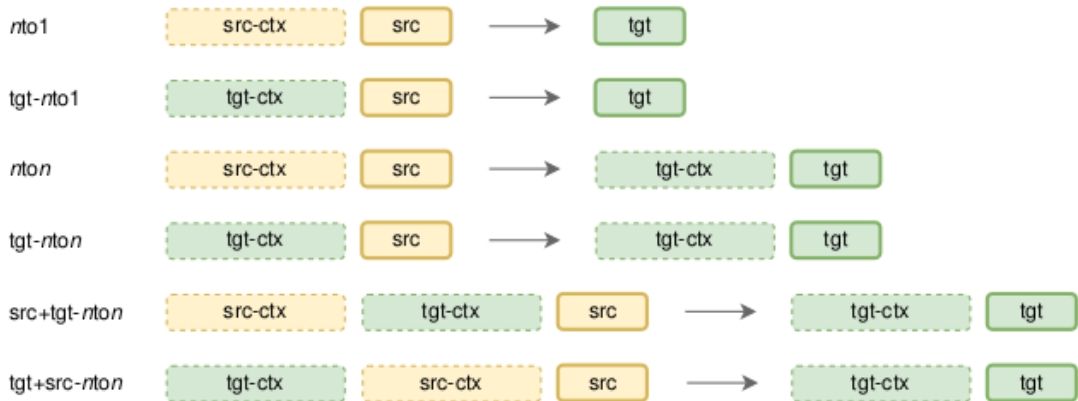


Figure 1: Schematic representation of a training instance for the different models. The yellow blocks represent the source language and the green blocks the target language. The dashed lines indicate context sentences; the continuous lines indicate the current sentence and its translation.

A Models Overview

To clarify the differences between model variants, Figure 1 provides a schematic view of the composition of a training instance for each type of concatenation-based model. We show the main building blocks and their ordering for both source and target sides.

B Machine-translated Target Context on Challenge Tests

To complement the results in Section 8, we evaluated the models on the challenge test sets using the machine-translated context instead of the reference translation in the test. Although this would be the process at inference time, as previously noted the challenge test sets depend on pre-established translation choices, in particular for coherence. A machine-translated context sentence might be entirely correct but differ from the specific translation choice the test has been designed for. The reference target context is thus typically provided as is on these test sets for standard approaches such as the *nton* model and we followed this protocol for our main results.

With these caveats in mind, we computed results in terms of BLEU and accuracy using machine translated-context on a subset of the challenge tests, with the results shown in Table 14. For this evaluation, we discarded the tests where the disambiguating information is present only in the target context, as this would lead to erroneous results, for the reasons mentioned above. Thus, the evalua-

tion was restricted to the GDR-SRC+TGT test for Basque-Spanish, and on the ellipsis-related tests for English-Russian. Although the contrastive results on these challenge tests might still be impacted by differing translation choices, the source context might contain sufficient information to compensate for these variations.

Using MT output impacted all the models that promoted the target context, in terms of both BLEU and accuracy scores, except in Basque-Spanish on BLEU where the loss was not statistically significant. However, these variants still outperformed the sentence-level baselines in a significant way across the board.

In English-Russian, the *src+tgt-nton* model using machine-translated context achieved better results than all other models on *Ellipsis VP*, excepting the same variant using reference translations. It was notably better than the *nton* and the *tgt-nton* models with reference target context. The situation is reversed on *Ellipsis infl.*, with significant losses for the *src+tgt-nton* (MT) model compared to *src+tgt-nton* (RF), and the *nton* model achieving better results with MT context. Note that the *nton* model also incurred significant losses in terms of accuracy when using MT context in this case. This is not unexpected, as the decoding process involves the target context in these models, with cascading divergences between the machine-translated target context and the expected context in the contrastive test. Note that this type of model is not impacted by the use of MT output in terms

	EU-ES		EN-RU			
	GDR-SRC+TGT		Ellipsis infl.		Ellipsis VP	
	BLEU	ACC.	BLEU	ACC.	BLEU	ACC.
Sentence-level	36.28	53.67	30.81	51.80	22.20	27.80
<i>nton</i> (RF)	40.45	77.67	36.97	75.20	29.59	62.60
<i>nton</i> (MT)	40.45	74.33	36.97	67.40	29.59	63.20
<i>tgt-nton</i> (RF)	39.05	72.67	40.69	70.00	30.75	60.00
<i>tgt-nton</i> (MT)	37.45	69.33	34.44	62.40	30.18	55.20
<i>src+tgt-nton</i> (RF)	41.25	78.67	40.98	77.20	35.84	77.60
<i>src+tgt-nton</i> (MT)	39.63	73.33	36.40	62.20	33.36	71.40

Table 14: Results on contrastive tests using reference (RF) and machine-translated (MT) context.

of BLEU, however, as the translated context is discarded after translation in non-contrastive evaluations.

In Basque-Spanish, the slight loss in BLEU between *src+tgt-nton* (RF) and *src+tgt-nton* (MT) was not statistically significant. In terms of accuracy, the losses were notable between these two models however, at over 5 points, but marginal between the *src+tgt-nton* (MT) and the *nton* (MT) models (1 point).

As previously discussed, contrastive tests are meant for a specific context, and evaluations with machine-translated output are only tentative. Different evaluation protocols would be needed to evaluate the use of MT context in a more principled and robust manner.

A Human Perspective on GPT-4 Translations: Analysing Faroese to English News and Blog Text Translations

Annika Simonsen and Hafsteinn Einarsson

University of Iceland

Sæmundargata 2, 102 Reykjavík, Iceland

{annika, hafsteinne}@hi.is

Abstract

This study investigates the potential of Generative Pre-trained Transformer models, specifically GPT-4, to generate machine translation resources for the low-resource language, Faroese. Given the scarcity of high-quality, human-translated data for such languages, Large Language Models' capabilities to produce native-sounding text offer a practical solution. This approach is particularly valuable for generating paired translation examples where one is in natural, authentic Faroese as opposed to traditional approaches that went from English to Faroese, addressing a common limitation in such approaches. By creating such a synthetic parallel dataset and evaluating it through the Multidimensional Quality Metrics framework, this research assesses the translation quality offered by GPT-4. The findings reveal GPT-4's strengths in general translation tasks, while also highlighting its limitations in capturing cultural nuances.

1 Introduction

In the past decade, the field of Natural Language Processing (NLP) has seen a dramatic shift with the introduction of the attention mechanism and Transformer models, profoundly influencing the domain of Machine Translation (MT) (Bahdanau et al., 2014; Vaswani et al., 2017). One of the foremost challenges in MT is the scarcity of high-quality, human-translated data for low-resource

languages. However, Large Language Models (LLMs) such as the Generative Pre-trained Transformer (GPT) models may present a solution to this challenge. These models are trained on vast amounts of data and have an impressive ability to generate native-sounding text, which they do by adapting based on the context presented in their training material (*in-context learning*) (Brown et al., 2020). Transformer models, such as GPT, are trained on multilingual data and have zero-shot translation capabilities which enables them to translate low-resource languages as well as high-resource languages. Therefore, this shift towards in-context learning signifies a breakthrough in NLP, where human-quality translation pairs can be generated without the input of a human translator. This has the potential of lowering the cost of making such data and improving the scalability of such an operation, which is vital for making smaller and more cost-effective models for MT.

This shift is particularly evident in the realm of MT datasets, where the gap between high-resource and low-resource languages remains a critical challenge. In the past, common methods to synthesize data for MT datasets were based on backtranslation (Sennrich et al., 2016; Poncelas et al., 2018; Poncelas et al., 2019). However, the quality of GPT models indicate that it is a better choice for synthesizing MT datasets (Hendy et al., 2023; Lyu et al., 2023).

The release of GPT-4 in March 2023 marked a significant milestone, with Jiao et al. (2023)'s pilot study that demonstrated its enhanced translation abilities in languages including English, German, Romanian, and Chinese, where its MT performance was comparable with state-of-the-art Neural Machine Translation (NMT) models. These results served as a motivation to explore the poten-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

tial of building MT datasets for low-resource languages with GPT models. Unlike traditional translation software, GPT-4 was not trained with the explicit purpose of MT, and in fact, it is not clear to what extent it was trained in MT at all. Further research, like Yang and Nicolai’s (2023), demonstrated the potential of applying GPT-generated synthetic data in the context of MT. Their models translated from German (a high-resource language) and Galician (a low-resource language). Their findings revealed that while models trained solely on natural data outperformed those trained solely on synthetic data, the best performing model was the one trained on a combination of both datasets. These findings encourage the validation of translation quality in GPT models for low-resourced languages such as Faroese.

The population of the Faroe Islands is approximately 54,500 people (Statistics Faroe Islands, 2024), with the large majority speaking Faroese as their L1. At this time, Faroese MT resources are lacking, and the existing resources are not sufficient for training high performing MT models (Simonsen et al., 2022). Focusing on Faroese, this paper explores GPT-4’s effectiveness in translating from Faroese to English. The potential of the GPT models raises the pivotal question: can the creation of MT resources for the low-resource language, Faroese, be automated, specifically through the capabilities of advanced models like GPT-4? To investigate this, the following contribution is made:

- A synthetic parallel dataset of 5,408 Faroese to English sentence pairs translated by GPT-4^{1,2}.
- A sample of 850 Faroese to English sentence pairs human evaluated and annotated with error labels from the Multidimensional Quality Metrics (MQM) framework by a native speaker^{3,4}.

¹https://huggingface.co/datasets/AnnikaSimonsen/GPT-4_FO-EN_parallel_news_sentences

²https://huggingface.co/datasets/AnnikaSimonsen/GPT-4_FO-EN_parallel_blog_sentences

³https://huggingface.co/datasets/AnnikaSimonsen/GPT-4_FO-EN_parallel_news_sentences_MQM

⁴https://huggingface.co/datasets/AnnikaSimonsen/GPT-4_FO-EN_parallel_blog_sentences_MQM

- An in-depth analysis of the types of reoccurring errors that GPT-4 makes when translating from Faroese to English.

These contributions provide a detailed human examination of GPT-4’s proficiency in translating Faroese to English, expanding on the current understanding of GPT-4’s translation capabilities of low-resource languages.

2 Previous work

2.1 Faroese Parallel Datasets

There has been some preliminary work done in Faroese MT, specifically within the domain of creating parallel training data. However, the state-of-the-art neural network MT models of today need vast amounts of training data, an obstacle that Faroese is still facing. The largest available parallel training data for Faroese is *Sprotin’s parallel corpus*⁵ which was published on GitHub in 2020 and contains over 100k sentences human-translated from English to Faroese. This initiative was part of an effort to encourage Google to include Faroese in the Google Translate application (Hvidfeldt, 2020). In response, Microsoft released a model trained on this dataset in their MT system called *Microsoft Translator*. An Icelandic NLP company, Miðeind, also released a Faroese MT model trained on the same data around the same time on their MT system called *Vélpýðing* (Símonarson et al., 2021). For both systems it was apparent that the model performance was not high, which was likely due to the small amount of training data. Faroese was never added to Google Translate and is also currently no longer supported on Vélpýðing. More recently, Meta launched the *No Language Left Behind* (NLLB) project which aims to bridge the gap in the performance between high- and low-resource languages in MT (Team et al., 2022). They published a series of open-sourced MT models called NLLB⁶ and a human-translated parallel dataset called FLORES-200⁷ which covers over 200 languages, including Faroese. The NLLB model’s capability to translate Faroese appears promising based on preliminary experiments made

⁵https://raw.githubusercontent.com/Sprotin/translations/main/sentences_en-fo.strict.csv

⁶<https://huggingface.co/facebook/nllb-200-distilled-1.3B>

⁷<https://huggingface.co/datasets/facebook/flores>

by the authors of this paper, a notable achievement considering the majority of its training data for Faroese is not genuinely parallel but is instead incorrectly aligned Faroese to English data.

Building upon the foundation of leveraging linguistic relations for enhancing machine translation in low-resource languages, recent studies have begun exploring the potential of utilizing phylogenetic information from high-resource languages within the same language family. For example, Snæbjarnarson et al. (2023) demonstrated that by incorporating resources from closely related Scandinavian languages, the performance of NLP tasks in Faroese could be substantially improved. This method marks a departure from the traditional 'one-size-fits-all' approach taken by widely used multilingual transformers like mBERT or XLM-R, advocating instead for a tailored strategy that considers the unique linguistic heritage of each language family. Such insights reveal the advantages of a more focused approach in data augmentation and model training, particularly for languages like Faroese that have fewer resources. Additionally, in Scalvini and Debess' (2024) upcoming publication, they highlight the effectiveness of GPT-SW3, a Scandinavian-focused LLM, in leveraging the linguistic similarities between Faroese and its Nordic counterparts to enhance translation accuracy and facilitate data augmentation.

2.2 Generating synthetic parallel data using GPT models

As mentioned in the introduction, there have been recent studies that explored using generative LLMs like ChatGPT or GPT-4 to create synthetic parallel data for training MT models. Inspired by findings that GPT-4 could match the translation abilities of commercial NMT systems, Yang and Nicolai (2023) explored training translation models for German and Galician using ChatGPT-generated synthetic data. In their study, they compared models trained on natural data from TED Talks with those trained on synthetic data, created by translating seed words and sentences into English via ChatGPT. Although models trained on real data performed better than those trained on synthetic data, those trained on a mix of real and synthetic data (augmented model) showed improved translation quality for both languages. Interestingly, synthetic Galician data yielded better translation quality than the German synthetic data. How-

ever, the study showed that there was a lower linguistic diversity in the synthetic data compared to natural data, evidenced by a lower type-token ratio (TTR), indicating a repetition of sentences and limited vocabulary use. This highlights challenges in leveraging LLMs for low-resource synthetic data creation. Nonetheless, supplementing real data with synthetic data remains a promising strategy for training MT models for low-resource languages (Poncelas et al., 2018; Poncelas et al., 2019).

```
function_descriptions = [
  {
    "name": "translation_analysis",
    "description": "The function analysis text that has been translated from Faroese to English. The input translation should be of exceptionally high quality.",
    "parameters": {
      "type": "object",
      "properties": {
        "sentence_analysis_list": {
          "type": "array",
          "items": {
            "type": "object",
            "properties": {
              "original": {
                "type": "string"
              },
              "translation": {
                "type": "string"
              }
            }
          }
        }
      }
    },
    "required": ["sentence_analysis_list"]
  }
]
```

Listing 1: JSON schema for translation analysis.

3 Experimental setup

This section outlines the methodology used to examine GPT-4's effectiveness in generating parallel data for Faroese, a language with limited resources. Firstly, the experiment involves generating synthetic parallel data using GPT-4. This output was then evaluated using the Multidimensional Quality Metrics (MQM) framework with a single native speaker as an annotator.

3.1 Prompting Approach

To generate the synthetic parallel data, a structured prompting approach with GPT-4 was employed, setting the temperature parameter to 0 to guarantee uniform and deterministic output. The experiment extracted information from Faroese news and blog texts through OpenAI's API, organizing it according to a specific JSON format (as detailed in Listing 1). This format instructed GPT-4 to translate texts sentence by sentence.

3.2 Data Preparation

3.2.1 GPT-4 Parallel Sentences

The parallel sentences generated by GPT-4 were derived from the Basic Language Resource Kit for Faroese 1.0 text corpus (Simonsen et al., 2022). During the first round, news texts from the online newspapers *Dimmalætting* and *Portalurin* were processed, with GPT-4 translating each document sentence by sentence, yielding 3,735 Faroese to English news sentence pairs. Subsequently, blog texts were translated using the same procedure, including works from *Egið Rúm* by Marna Jacobsen⁸, *BAVS* by Bergljót av Skarði⁹, and *BirkBlog* by Birgir Kruse¹⁰, resulting in 1,673 sentence pairs from blogs. The aim was to capture a diverse representation of GPT-4’s translation skills by combining news and blog texts, acknowledging their distinct genres. In total, there were 5,408 generated sentence pairs.

3.3 Human Evaluation

A subset of the GPT-4 generated parallel data was sampled for human evaluation using the MQM framework¹¹. MQM incorporates more than twenty traditional translation quality metrics and provides a detailed catalogue of over 100 potential issues for assessing translations and source texts. It is designed as a flexible master list from which specific issues can be chosen based on the translation task at hand, allowing for customization to meet diverse requirements. In the context of this study’s MT evaluation, a tailored version of the MQM framework, as adapted by Freitag et al. (2021) was employed. An overview of the MQM error categories utilized by Freitag et al. (2021), along with their descriptions, is presented in Table 1.

The sample that was chosen for human evaluation was created by choosing articles randomly and then evaluating the chosen articles, sentence by sentence. There was only one annotator, author of this paper, who is a linguist and native speaker of Faroese. In total, 425 news sentence pairs and

425 blog sentences pairs were human evaluated. The evaluation was carried out in a Google Sheet spreadsheet (see Figure 1). To calculate the MQM score, the official MQM spreadsheet was used. This spreadsheet contains all relevant formulas to calculate the MQM score, also known as the *Overall Quality Score* (OQS).

4 Results

The results for the MQM evaluation is summarized in Table 2. Overall, the quality of translations is high as indicated by the MQM score or *Overall Quality Score*.

As seen in Table 2, the predominant severity level assigned for MQM was *minor*. This classification was used when translations were not technically accurate but still conveyed the intended meaning. The *major* category was designated for errors in translation that obscured or altered the meaning, while *critical* was reserved for when the translation got offensive or dangerously misinformed. Notably, major accuracy errors occurred more frequently in blogs than in news articles, often arising in idiomatic expressions and set phrases. In the case of news articles, there were three instances classified as *critical*¹²:

1) *Example sentence containing critical error from Portalurin article.*

- **FO:** *Somuleiðis skulu dagfóringar gerast á Vágs høll við máling og wc til røðslutarna skal gerast í ganginum millum VB húsið og Vágs Høll.*
- **ENG:** "Similarly, updates should be made to the Vágur hall with regards to painting and a toilet for **disabled** that should be made in the corridor between the VB house and Vágur Hall."
- **GPT-4:** "Similarly, updates should be made to Vágur hall with painting and a toilet for **the pipe players** should be made in the corridor between the VB house and the Vágur Hall."

The original sentence contains a spelling mistake; *røðslutarna* is supposed to be spelled *rørslutarnað* which translates to "disabled". GPT-4 translated this term to "the pipe players".

⁸Jacobsen shares insights from her personal life, coupled with reviews of music, books and movies. Available at: <https://marnakj.wordpress.com/>.

⁹BAVS is centered on personal experiences, culture, and travel. Available at <https://b-av-s.blogspot.com/>

¹⁰A blog focusing on cultural events along with reviews of movies, music, and more. Available at: <https://birkblog.blogspot.com/>

¹¹<https://www.qt21.eu/>

¹²The order in the sentence examples is as follows: **FO** (original Faroese sentence from dataset), **ENG** (a translation provided by an author of this paper) and **GPT-4** (GPT-4’s translation of the original Faroese sentence).

Error Category	Description
Accuracy	
Addition	Translation includes information not present in the source.
Omission	Translation is missing content from the source.
Mistranslation	Translation does not accurately represent the source.
Untranslated text	Source text has been left untranslated.
Fluency	
Punctuation	Incorrect punctuation (for locale or style).
Spelling	Incorrect spelling or capitalization.
Grammar	Problems with grammar, other than orthography.
Register	Wrong grammatical register.
Inconsistency	Internal inconsistency (not related to terminology).
Character encoding	Characters are garbled due to incorrect encoding.
Terminology	
Inappropriate for context	Terminology is non-standard or does not fit context.
Inconsistent use	Terminology is used inconsistently.
Style	
Awkward	Translation has stylistic problems.
Other	Any other issues.
Source error	An error in the source text.
Non-translation	Impossible to reliably characterize the 5 most severe errors.

Table 1: Overview of MQM label hierarchy.

Faroese	English translation	Errors	Severity
Danska rokktrioin, sum legði ríkið fyrri sínar fætur í 90'unum, eigur eitt heilt serligt pláss í hjartanum á mongum føroyingi.	The Danish rock trio, which conquered the country in the 90s, holds a special place in the hearts of many Faroese.	style(føroyingi-Faroese/Faroe Islanders)	minor
Trioin gav út plátuna Dizzy Mizz Lizzy í 1994, og hon streyk upp á tónleikatindarnar alt fyrri eitt.	The trio released the album Dizzy Mizz Lizzy in 1994, and it shot up the music charts immediately.		
Útgávan var serstøk, tí hon var á tremur við hittum.	The release was special because it was on par with others.	accuracy(á tremur við hittum-on par with others/full of hits)	major

Figure 1: Figure showing the human evaluation method for the GPT-4 generated parallel sentences in Google Sheets.

2) *Example sentence containing critical error from Dimmalætting article.*

- **FO:** *Orsøkin er, at svenski Umhvørvisflokkurin, ið var í stjórn saman við Sosialdemokratunum, hevur vent samgonguni bakið, eftir at tað gjørdist greitt, at borgarliga andstøðan fekk ein meiriluta við at atkvøða sína fíggarlóg ígjøgnum.*
- **ENG:** "The reason is that the Swedish Environmental Party, which was in government with the Social Democrats, has turned its back on the coalition, after it became clear that **the civil opposition** got a majority by voting their budget through."
- **GPT-4:** "The reason is that the Swedish Environmental Party, which was in government with the Social Democrats, has turned its back on the coalition, after it became clear that **the bourgeois opposition** got a majority by voting their budget through."

The word for "civil opposition" has been translated into "bourgeois opposition", which could have negative connotations. The third example is in the same article where "civil budget" was translated into "bourgeois budget".

In the blog texts, a single critical error was identified in the sample, involving the mistranslation of

the term *at ræsa* — the traditional Faroese method of fermenting meat through dry-aging. It was incorrectly translated as "raw". This misinterpretation could be seen as dangerously misleading and potentially harmful to someone's health, especially if the food had not been pre-cooked as could be inferred from the context given correct world-knowledge:

3) *Example sentence containing critical error from BirkBlog.*

- **FO:** *Eg vildi smakka ræstu pylsuna.*
- **ENG:** "I wanted to taste the **Faroese dry-aged** sausage."
- **GPT-4:** "I wanted to taste the **raw** hot dog."

The Overall Quality Score (OQS), as detailed in Table 2, serves as a metric for assessing translation quality. It is derived through a systematic procedure: annotators input error annotations into a matrix (see Figure 2), assigning them numerical values based on error type and severity, to obtain the Absolute Penalty Total (ABT). The OQS calculation incorporates several factors, including the Per-word Penalty Total, calculated by dividing the ABT by the total word count (EWC); the Overall Normed Penalty Total, which adjusts the per-word penalty in relation to the total number of reference

Category	News Sentences (425)			Blog Sentences (425)		
	Minor	Major	Critical	Minor	Major	Critical
Accuracy	43	26	1	31	54	1
Fluency	9	0	0	8	5	0
Terminology	76	14	2	13	8	0
Style	35	2	0	11	3	0
MQM Score	94.41			88.38		

Table 2: MQM evaluation results for 425 news sentences and 425 blog sentences. The MQM score is also known as the Overall Quality Score (OQS). The weights are minor (-1), major (-5) and critical (-25). A higher score (with a maximum of 100) corresponds to better performance.

	A	B	C	D	E	F	G	H
1	MQM Scorecard: Top-Level Error Typology with 4 Severity Levels							
2								
3			Error Severity Levels:	Neutral	Minor	Major	Critical	Error Type Penalty Total
4			Severity Penalty Multipliers:	0	1	5	25	
5	ET Nos	Error Types	Error Counts				ET Weights	ETPTs
6	1	Terminology	2	7	7	0	1.0	42.0
7	2	Accuracy	4	14	7	1	1.0	74.0
8	3	Linguistic conventions	1	23	9	0	1.0	68.0
9	4	Style	5	7	3	0	1.0	22.0
10	5	Locale convention	1	12	5	0	1.0	37.0
11	6	Audience appropriateness	0	2	1	0	1.0	7.0
12	7	Design and markup	0	6	1	0	1.0	11.0
13	8	Custom						
14							Absolute Penalty Total:	261.00
15								
16		Evaluation Word Count:	10184				Per-Word Penalty Total:	0.0256
17		Reference Word Count:	1000				Overall Normed Penalty Total:	25.63
18		Scaling Parameter (SP):	1.00				Overall Quality Score:	97.44
19		Max. Score Value:	100.00					
20		Threshold Value:	85.00				Pass/Fail Rating:	Pass

Figure 2: Figure showing the Overall Quality Score card from <https://themqm.org/>.

words; and the Overall Quality Fraction, achieved by dividing the ABT by the EWC. The final Overall Quality Score is computed by subtracting the result of multiplying the per-word penalty score by the highest possible score from 1, thereby converting the score into a more recognizable percentage format. This approach integrates a meticulous evaluation of translation inaccuracies with a comprehensive scoring framework to measure and express the quality of translations quantitatively.

In Table 2, the Overall Quality Score demonstrates high performance in translation quality for both text genres, with news articles achieving a score of 94.41/100 and blogs receiving 88.38/100. According to the MQM framework, a score within the range of $94 \leq x < 98$ signifies a high level of quality, whereas scores in the range of $80 \leq x < 94$ are indicative of a good quality level, as outlined in Talhadas (2023). Following this, a qualitative analysis is provided to examine the specific types of translation errors GPT-4 made while translating from Faroese to English.

4.1 Most Common FO–EN Translation Errors by GPT-4

There is a general pattern in the types of errors that GPT-4 makes when translating from Faroese to English. A prominent error involves the translation of the Faroese term for the Danish currency used in the Faroe Islands, *krónur*. GPT-4 often renders these as "crowns" or "kroner", whereas the conventional translation should be "DKK" or "Danish crowns." Additionally, we observed four times that *føroyingur* was translated to "the Faroese" in the sample, but a more precise translation would be "a Faroe Islander" or "a Faroese person." A check on the rest of the translated data revealed that this was a common mistranslation.

Another consistent error is the translation of *korona* to "corona." While not incorrect, "COVID-19" is the term more frequently used in English news articles, making it a more suitable translation in those contexts. Given that "COVID" was not adopted into Faroese during the pandemic, Faroese news texts use *korona* instead.

Subsequent sections will detail the other preva-

lent errors, categorized by their types, to provide a comprehensive overview of the translation challenges encountered. See Table 3 in the Appendix for a quantitative analysis of the errors analyzed in this section.

4.1.1 Named Entities (NEs)

GPT-4 did not consistently translate all NEs incorrectly, but did manage to correctly translate certain NEs, such as names of institutions, e.g. *Ráðið fyri Ferðslutrygd* ("Council for Traffic Safety"). It also frequently accurately converted people's names from the dative to the nominative case in English, such as translating *Mariu* to "Maria". Nevertheless, there are many examples of translation errors with NEs. For instance, the short form name *Setrið* for *Fróðskaparsetrið* was incorrectly translated literally as "Center" rather than "The University of the Faroe Islands" or simply "Setrið". Additionally, *Okkara Voxbotn* was inaccurately translated to "Our Voxbotn" instead of preserving its original name, which is associated with a brewery named *Okkara* that sponsors a music festival that takes place in the harbour named *Voksbotn*. These examples illustrate the nuanced difficulties GPT-4 faces with NEs in the context of Faroese to English translation.

4.1.2 Correct Translation, Wrong Terminology

While GPT-4 frequently chooses accurate translations for words, it often selects the wrong terms. For instance, in some contexts *skeið* is translated as "course" when it is supposed to be "workshop," and *eldraøki* is translated as "elderly area" instead of "elderly affairs". The term *øki* can be translated as "area" only when it is referring to a physical place. In this context, the term was used in a sentence from an article about a financial budget of a town, where the taxes had been increased to cover elderly affairs. Therefore, while these translations are technically correct, they are not entirely appropriate in these specific contexts.

4.1.3 Idioms and Fixed Phrases

GPT-4 often encounters difficulties with idiomatic expressions and fixed phrases, particularly in the Faroese blog texts. These phrases frequently undergo literal translation, which misses their nuanced meanings. For instance, the phrase *at fáa sær okkurt gott* is directly translated as "getting oneself something good" rather than capturing the

intended meaning of "getting something to eat." Similarly, a well-known Faroese phrase, *er ikki sum at siga tað*, intended to convey that something is not easy, is translated by GPT-4 in a literal manner as "is not like saying it."

4.2 Icelandicisms

GPT-4 often confounds Faroese with Icelandic, likely due to the fact that GPT-4 has been trained on significantly more Icelandic data than Faroese. This leads to what we term "Icelandicism", where translations mistakenly apply Icelandic meanings. Examples include translating *menning* as "culture" instead of "progress", *sætti* as "sweet" instead of "sixth" and *bleytur* as "wet" rather than "soft". Here, it is presumed that the Faroese word *sætti* is confused with the Icelandic word *sætur* and the Faroese word *bleytur* is conflated with the Icelandic word *blautur*.

4.2.1 Cultural Context

GPT-4 often misses the cultural nuances in its translations, leading to misunderstandings of certain terms. For example, it interprets *ríkið* as "country" instead of "kingdom". In Faroese contexts, *ríkið* typically refers to the Danish Kingdom rather than the Faroe Islands. This misinterpretation might also reflect an Icelandicism. Other Faroese terms that are commonly mistranslated include *fiskaplasið*, which refers to a stone-paved area for drying fish but gets translated as "fish place," and *hoyggjhús*, which is translated as "living room" instead of "hay barn". Occasionally, these culturally specific terms are translated into nonsensical words. For instance, the Faroese term for "paternal granduncle", *abbabeiggi*, was erroneously translated as "abbess." The term *abbabeiggi* is culturally significant, as, although Icelandic also has a term for the brother of your grandfather, *afabróður*, it is not as commonly used as in Faroese. Notably, Danish lacks a term for this specific type of granduncle. Another example of a mistranslation is the Faroese word for a national dish, pilot whale steak (*grindabúffur*), which was translated to the nonsense word "grindabuffi".

This examination underscores that although GPT-4's FO-EN translations are of commendable quality, they exhibit specific and frequently recurring mistakes, notably in handling cultural subtleties and idiomatic expressions. Further qualita-

tive analyses of errors are deferred to the appendix.

5 Discussion

According to the human evaluation of the GPT-4 translation data, GPT-4 has demonstrated its proficiency in translating from Faroese to English, especially in the context of news articles. However, the translations are not perfect and we address the limitations later in this section. This finding of translation quality is consistent with recent research indicating GPT-4’s effectiveness in translating from low-resource languages to English, as highlighted in studies by Bang et al. (2023), Jiao et al. (2023), and Yang and Nicolai (2023). The model’s particular strength in news translation is likely due to its extensive training on a wide array of news texts, which is abundantly available online for collection. However, when it comes to blog texts, which are rich in idiomatic expressions and fixed phrases, GPT-4’s performance dips slightly (from 94.41 to 88.38). This drop in performance suggests that while GPT-4 can generate high-quality translations, its capability diminishes with content that heavily features language-specific idioms and cultural nuances. Yet, the synthetic parallel sentences generated by GPT-4 present a "Silver Standard" resource for training MT models for Faroese, complementing the "Gold Standard" human-translated data. Although not a substitute for human translation, the combination of synthetic and human-generated data could potentially enhance the training materials available for Faroese MT models (for German and Galician, see Yang and Nicolai (2023); for English, German and Turkish, see Sennrich et al. (2016). However, to fully assess the impact of GPT-4 generated parallel data on MT model performance, a larger dataset would be ideal. For this study, the collection was limited to 5,408 sentence pairs due to cost considerations and the licensing restrictions imposed by OpenAI on their model’s output¹³.

During the study, a preliminary experiment was conducted to see how well GPT-4’s translation performed from English into Faroese, which revealed significant limitations. The model often failed to construct grammatically correct Faroese sentences, frequently producing outputs that appeared to be an amalgamation of Icelandic and Faroese. This finding corroborates previous re-

¹³At the time of usage, gpt-4-0613 cost \$0.03 for every input token and \$0.06 for every output token.

search indicating that GPT-4’s capabilities in translating from English to low-resource languages remain constrained. Studies by Hendy et al. (2023), Lyu et al. (2023), Jiao et al. (2023), and Yang and Nicolai (2023) have similarly documented these challenges, reinforcing the observation that GPT-4’s performance in such translation tasks is not yet satisfactory.

6 Limitations

6.1 GPT-4’s Splitting of Sentences

GPT-4 did not consistently split the Faroese text into sentences although it was explicitly instructed to do so using our function-callin approach to extract output in a structured manner. An analysis of a random selection of 500 rows from the parallel dataset generated by GPT-4 revealed 29 cases of improper sentence division. This means that GPT-4 incorrectly split the Faroese sentences 5.8% of the time. This could also be related to the reason why GPT-4 struggled with translating NEs. NEs are notoriously difficult for MT models to handle accurately, and Named Entity Recognizers (NERs) are often employed alongside these models to enhance performance (Babych and Hartley, 2003). In the early stages of developing the GPT-4 parallel data for this experiment, attempts were made to have GPT-4 label the Faroese text with NE labels. However, these efforts were unsuccessful, leading to the exclusion of this step from the process. This difficulty likely stems from GPT-4’s inadequate ability to recognize Faroese NEs, contributing to its struggles with their translation.

6.2 Systematic Translation Errors

While GPT-4 delivers translations of high quality from Faroese to English, it is worrying to see that the errors it makes are often specific to Faroese context and culture. These mistakes do not seem to be random but show a pattern that could negatively affect the efficacy of an MT model trained with such data. A possible remedy might have been to enrich GPT-4’s contextual understanding, perhaps by feeding it Wikipedia articles that encapsulate key facts about Faroese culture and the Faroe Islands, or by exposing it to Faroese texts across different genres. This enhanced prompting strategy, not dissimilar to few-shot prompting (Brown et al., 2020) could have helped GPT-4 in situations where its grasp of context and global knowledge fell short. Ultimately, the synthetic parallel data

produced by GPT-4 ought to be considered as "Silver Standard", rather than "Gold Standard" data, which is typically human-translated. Drawing parallels to the findings of Yang and Nicolai (2023) regarding German and Galician ChatGPT-generated parallel data, it becomes apparent that Silver Standard Data holds value, particularly when combined with human-translated data for training Faroese MT models to maximize performance.

6.3 Lack of MQM annotators

It is crucial to acknowledge that since there was only one annotator, the MQM scores should not be compared with those from projects that had multiple annotators, under the assumption that the result for Faroese is not as robust as for other languages (i.e. due to potential individual biases). The primary aim of conducting the MQM evaluation was to delve into error analysis and obtain a comprehensive understanding of the translation quality. Additionally, the Faroese annotator chose not to apply the "neutral" error weight during the MQM assessment, a deviation from conventional practices. This decision was made because labeling an error as "neutral" seemed inappropriate when such a categorization is typically reserved for instances deemed not to be the translator's fault and, in this case, the translator is a language model. Looking back, this neutral category might have been applicable for source text errors, such as typos, but ultimately, these errors were given the label "source error", so the resulting score was not affected as "source errors" count the same as a "neutral" error. Only 12 source errors were found in total, and only four of them resulted in a translation error. Determining whether an error is attributable to the language model presents its own challenges. Furthermore, given that the MQM framework is designed for evaluating both human- and machine translation, applying it uniformly to both can be problematic. For instance, human translators often tailor their work to a client's specific style requirements, which can range from general and succinct to verbatim translations, depending on whether clarity or fidelity is prioritized. Current MT models, however, lack the capability to adjust their output based on stylistic preferences without specific training for each requirement. Nevertheless, LLMs like GPT-4 have shown the ability to adapt to given instructions, suggesting they can be directed to follow certain styles or formalities. Future research

may need to reconsider how we evaluate LLMs like GPT-4, taking into account their unique capabilities and limitations.

6.4 OpenAI and Model Ownership

OpenAI's terms of use (OpenAI, 2023) stipulate that while users are granted ownership of the output generated by its services, there are restrictions when it comes to using GPT-4 output for model training. Specifically, users are prohibited from using the output to develop models that compete with OpenAI. As a result, the authors of this paper limited their generation to approximately five thousand parallel sentences with GPT-4, as they would not have been able to share any models fine-tuned using this training data. Similarly, AI META's Llama 2 model permits derivative works but forbids their use in enhancing language models other than Llama (Meta AI, 2023), which would have prevented the authors from sharing their fine-tuned MT models had they used Llama 2 instead of GPT-4.

However, there are open-source LLMs that serve as alternatives, some of which aim to address the lack of diversity in the text used to train LLMs. Notable efforts include AI Sweden's GPT-SW3 (Swedish Government, 2023), focused on Nordic languages, and the upcoming Horizon Europe funded TrustLLM, aiming for an open, trustworthy, and Germanic language-focused LLM¹⁴. AI Sweden offers a flexible license for GPT-SW3 (AI Sweden, 2023), exemplifying the push towards democratizing LLM access. It is worth noting that there are significant differences in parameter size between these models, with GPT-SW3's largest instruct model having 20B parameters, Llama 2's biggest instruct model having 70B parameters, and GPT-4 believed to have over a trillion parameters. Another recently published open model is Mistral's Mixtral 8x7B, which is under the Apache 2.0 license and is reported to either match or outperform Llama 2 and GPT-3.5 on most standard benchmarks (Mistral AI team, 2023). These open models provide potential alternatives for future work in automating Faroese NLP resource creation.

6.5 Future Work

This research has identified several promising directions for future work in Faroese MT. Firstly,

¹⁴<https://trustllm.eu/>

the potential of synthetic parallel data produced by LLMs like GPT-4 for Faroese remains largely unexplored. Future efforts should focus on creating a larger corpus of synthetic parallel sentences covering a wider range of text genres beyond news and blogs. This approach would provide insights into how effectively such data can train more robust MT systems. However, licensing restrictions associated with some LLMs may necessitate a shift towards openly available models, such as GPT-SW3, the forthcoming Germanic LLM from the TrustLLM project, or Mistral’s Mixtral 8x7B. These open models would facilitate the generation of larger datasets and ensure the ability to freely share and distribute the resulting works, aligning with research efforts aimed at enhancing NLP capabilities for low-resource languages like Faroese. Scalvini and Debess (2024) have demonstrated the merits of using language-family-specific models, such as GPT-Sw3, in refining translation accuracy and facilitating data augmentation efforts for Faroese.

Secondly, there is currently no human-translated parallel dataset for Faroese derived from monolingual Faroese texts. Existing datasets, such as FLORES-200 and the Sprotin parallel corpus, are translations from English and do not accurately reflect Faroese-specific expressions and terminologies. Consequently, the synthetic parallel data generated by GPT-4 also falls short in capturing these unique Faroese nuances. Therefore, developing a human-translated parallel dataset centered around Faroese monolingual content, with an emphasis on capturing the richness of Faroese cultural and linguistic elements, would be highly advantageous for future research in Faroese MT.

Finally, recent advancements in models like Gemini 1.5 Pro, which can process exceptionally long contexts, have opened up new prospects for MT in Faroese. Gemini 1.5 Pro has demonstrated its ability to learn new languages from a minimal set of instructional materials. Specifically, with only 500 pages of linguistic documentation and approximately 400 parallel sentences, it managed to learn and translate from English to Kalamang, a critically low-resource language with minimal online presence, achieving translation quality comparable to human learners (Gemini Team, 2024). This success suggests that for Faroese, leveraging Faroese grammar books and lexical resources in the translation context could make high-quality

translation not only feasible but also efficient. This promising approach warrants further investigation in future research.

7 Conclusion

In conclusion, this paper has demonstrated the potential of GPT models like GPT-4 in generating synthetic parallel data, potentially mitigating the scarcity of high-quality, human-translated datasets. Through a detailed analysis of GPT-4’s translation from Faroese to English, including a synthetic parallel dataset and an MQM framework-based evaluation, we have uncovered both strengths and limitations of employing GPT models for MT. While GPT-4 shows promise in generating translations that could serve as valuable training data, challenges remain, particularly with translations that involve cultural and contextual nuances. This exploration not only contributes to the understanding of GPT models’ capabilities in translating low-resource languages but also sets the stage for future research directions. By integrating synthetic and human-generated data, there’s potential to enhance MT models for Faroese, pushing the boundaries of accessibility and quality in MT for low-resource languages. This study underscores the necessity for ongoing research to fully leverage the capabilities of advanced models like GPT-4, aiming for a future where no language is left behind in the digital age.

Acknowledgments

AS was supported by the European Commission under grant agreement no. 101135671. We thank the reviewers for their constructive and helpful comments, which have significantly improved the quality and clarity of our manuscript.

References

- AI Sweden. 2023. AI Sweden’s LLM AI Model License Agreement for GPT-SW3. Accessed: December 2023.
- Babych, Bogdan and Anthony Hartley. 2003. Improving machine translation quality with automatic named entity recognition. In *Proceedings of the 7th International EAMT workshop on MT and other language technology tools, Improving MT through other language technology tools, Resource and tools for building MT at EACL 2003*, Columbus, Ohio.

- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Bang, Yejin, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A multitask, multilingual, multimodal evaluation of ChatGPT on reasoning, hallucination, and interactivity. In Park, Jong C., Yuki Arase, Baotian Hu, Wei Lu, Derry Wijaya, Ayu Purwarianti, and Adila Alfa Krisnadhi, editors, *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 675–718, Nusa Dua, Bali, November. Association for Computational Linguistics.
- Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Larochelle, H., M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Daems, Joke, Orphée De Clercq, and Lieve Macken. 2017. Translationese and post-editeese: How comparable is comparable quality? *Linguistica Antverpiensia New Series-Themes in Translation Studies*, 16:89–103.
- Freitag, Markus, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.
- Gemini Team. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. Technical report, Google DeepMind. Technical Report.
- Hendy, Amr, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.
- Hvidfeldt, Jón Brian. 2020. Faroese language to feature on Google Translate, 10. Accessed: December 2023.
- Jiao, Wenxiang, Wenxuan Wang, Jen-tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Is ChatGPT a Good Translator? Yes With GPT-4 As The Engine. *arXiv preprint arXiv:2301.08745*.
- Lyu, Chenyang, Jitao Xu, and Longyue Wang. 2023. New trends in machine translation using large language models: Case examples with ChatGPT. *arXiv preprint arXiv:2305.01181*.
- Meta AI. 2023. Llama 2 community license agreement. <https://ai.meta.com/llama/license/>. Accessed: December 2023.
- Mistral AI team. 2023. Mixtral of experts: A high quality sparse mixture-of-experts. <https://mistral.ai/news/mixtral-of-experts/>. Accessed: January 2024.
- OpenAI. 2023. Terms of use. <https://openai.com/policies/terms-of-use>. Accessed: December 2023.
- Poncelas, Alberto, Dimitar Shterionov, Andy Way, Gideon Maillette de Buy Wenniger, and Peyman Passban. 2018. Investigating backtranslation in neural machine translation. In Pérez-Ortiz, Juan Antonio, Felipe Sánchez-Martínez, Miquel Esplà-Gomis, Maja Popović, Celia Rico, André Martins, Joachim Van den Bogaert, and Mikel L. Forcada, editors, *Proceedings of the 21st Annual Conference of the European Association for Machine Translation*, pages 269–278, Alicante, Spain, May.
- Poncelas, Alberto, Maja Popović, Dimitar Shterionov, Gideon Maillette de Buy Wenniger, and Andy Way. 2019. Combining PBSMT and NMT back-translated data for efficient NMT. In Mitkov, Ruslan and Galia Angelova, editors, *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019)*, pages 922–931, Varna, Bulgaria, September. INCOMA Ltd.
- Scalvini, Barbara and Iben Nyholm Debess. 2024. Evaluating the potential of language-family-specific generative models for low resource data augmentation: a Faroese case study. In *Proceedings of The 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING '24)*, Torino.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In Erk, Katrin and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany, August. Association for Computational Linguistics.
- Símonarson, Haukur Barri, Vésteinn Snæbjarnarson, Pétur Orri Ragnarsson, Haukur Páll Jónsson, and Vilhjálmur Þorsteinsson. 2021. Miðeind’s WMT 2021 submission. *arXiv preprint arXiv:2109.07343*.

Simonsen, Annika, Sandra Saxov Lamhauge, Iben Nyholm Debess, and Peter Juel Henriksen. 2022. Creating a Basic Language Resource Kit for Faroese. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4637–4643.

Snæbjarnarson, Vésteinn, Annika Simonsen, Goran Glavaš, and Ivan Vulić. 2023. Transfer to a low-resource language via close relatives: The case study on Faroese. In Alumäe, Tanel and Mark Fishel, editors, *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 728–737, Tórshavn, Faroe Islands, May. University of Tartu Library.

Statistics Faroe Islands. 2024. Population. Accessed: February 2024.

Swedish Government. 2023. Regeringen tillsätter en AI-kommission för att stärka svensk konkurrenskraft, 12. Press release from Finansdepartementet, Statsrådsberedningen.

Talhadas, Paulo. 2023. Quality reports - monitor the quality results for your translations, 11. Accessed: December 2023.

Team, NLLB, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No Language Left Behind: Scaling Human-Centered Machine Translation.

Vanmassenhove, Eva, Dimitar Shterionov, and Matthew Gwilliam. 2021. Machine translationese: Effects of algorithmic bias on linguistic complexity in machine translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2203–2213.

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Yang, Wayne and Garrett Nicolai. 2023. Neural Machine Translation Data Generation and Augmentation using ChatGPT. *arXiv preprint arXiv:2307.05779*.

Zhang, Mike and Antonio Toral. 2019. The effect of translationese in machine translation test sets. In *Proceedings of the Fourth Conference on Machine*

Translation (Volume 1: Research Papers), pages 73–81, Florence, Italy, August. Association for Computational Linguistics.

A Quantitative Error Analysis

We count the most common errors and display them in Table 3. Furthermore, we refer to two additional categories of common errors found in GPT-4’s translations from Faroese to English.

Translationese

GPT-4 occasionally generates translations that come across as awkward or grammatically incorrect in English, an issue commonly encountered in MT referred to as "machine-translationese" (Zhang and Toral, 2019; Daems et al., 2017; Vanmassenhove et al., 2021). In the case of the Faroese to English translations, GPT-4 sometimes opts for a literal, word-for-word translation approach, leading to syntax that sounds unnatural. For example:

- **FO:** Illgruni er tó um, at sjeý onnur eisini eru smittað við nýggja frábrigðinum, **skrivar Ritzau.**
- **GPT4** "However, there is suspicion that seven others are also infected with the new variant, **writes Ritzau.**"

In this case, it is more natural to choose the word order, "Ritzau writes". However, it is worth to note that this type of error is possibly not thought of as an error by some, because it could in reality be a question of style-preference.

Inappropriate Register

Finally, GPT-4 sometimes opts for translations that carry an inappropriate tone, especially noticeable in formal settings such as news articles. For instance, *andaðist* is translated into the more colloquial "died" rather than the more fitting and respectful "passed away". This discrepancy in tone becomes particularly evident in news reporting, where a certain level of formality is anticipated. However, it is crucial to acknowledge that the register and genre of the text were not defined when prompting GPT-4.

Category	News Sentences (425)	Blog Sentences (425)
Correct-translation-wrong-terminology	89	20
NEs	26	5
Cultural context	16	18
Idioms and fixed phrases	11	19
DKK	19	0
Translationese	10	7
Icelandicism	6	7
Source error	8 (2)	2
Faroese	3	1
Inappropriate register	2	0
COVID	1	0
Other	41	56

Table 3: Detailed evaluation results for specific categories in translations of 425 news sentences and 425 blog sentences. The figures in parentheses indicate the count of translation errors within the total reported for that category.

ReSeTOX: Re-learning attention weights for toxicity mitigation in machine translation

Javier García Gilabert, Carlos Escolano

Universitat Politècnica de Catalunya
{javier.garcia.gilabert,
carlos.escolano}@upc.edu

Marta R. Costa-Jussà

FAIR, Meta
costajussa@meta.com

Abstract

Our proposed method, RESETOX (REdo SEArch if TOXic), addresses the issue of Neural Machine Translation (NMT) generating translation outputs that contain toxic words not present in the input. The objective is to mitigate the introduction of toxic language without the need for re-training. In the case of identified added toxicity during the inference process, RESETOX dynamically adjusts the key-value self-attention weights and re-evaluates the beam search hypotheses. Experimental results demonstrate that RESETOX achieves a remarkable 57% reduction in added toxicity while maintaining an average translation quality of 99.5% across 164 languages. Our code is available at: <https://github.com/mt-upc/ReSeTOX>

WARNING: the current paper contains examples that may be offensive.

1 Introduction

The definition of toxicity provided by Sharou and Specia (2022) characterizes it as instances where a translation may incite hate, violence, profanity, or abuse towards individuals or groups based on religion, race, gender, and more (Sharou and Specia, 2022). Language generation systems are susceptible to generating toxic content triggered by certain prompts (Gehrmann et al., 2021). Unlike Machine Translation (MT) systems that are conditioned on a given source input, unconditioned language generation systems are more susceptible to this safety

concern. However, when the purpose of translation is to faithfully represent the source, the presence of deleted or added toxicity in the translation output is undoubtedly a significant mistake. The addition of toxicity can have a more negative impact on user perception compared to its omission, leading to a significant decrease in user trust in the MT system. Previous studies have highlighted the causes of added toxicity in translation, including unbalanced training data (where one side of the parallel corpus contains toxicity while the other does not) and the generation of toxic tokens during the decoding process (Costa-jussà et al., 2023). Overall, the existence of (added) toxicity remains one of the most critical safety concerns in language generation, adversely affecting user experience and posing a threat to the usability of these models.

Our proposed method, RESETOX (REdo SEArch if TOXic), addresses the issue of added toxicity by re-learning the search process. Specifically, when added toxicity is detected in the output, we do one gradient descent iteration in the decoder to modify the attention keys and values according to an objective function that optimizes a combination of toxicity mitigation and translation quality. Then, we re-score the hypothesis from the beam search. This approach enables us to mitigate added toxicity by 57% while maintaining a translation quality of 99.5%. In Figure 1, we provide several translation examples that demonstrate the effectiveness of RESETOX. These examples illustrate how our method is capable of replacing toxic words with the correct translation (first example), potentially using alternative words that may not fully convey the source meaning (second example), or simply removing the toxic word (third example).

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

Source: I have a friend who is a stinky guy.
Baseline: J'ai un ami qui est un gars putain.
+RESETOX: J'ai un ami qui est un gars puant.
Ref: J'ai un ami qui est un gars puant.
Source: I have friends who are dreadlocked grandfathers.
Baseline: Tengo amigos que son abuelos con tetas.
+RESETOX: Tengo amigos que son abuelos atrevidos.
Ref: Tengo amigos que son abuelos con rastas.
Source: It's hard being a shaggy veteran.
Baseline: És difícil ser un veterà de merda.
+RESETOX: És difícil ser un veterà.
Ref: És difícil ser un veterà pelut.

Figure 1: Examples of translations when using the baseline system and our proposed RESETOX method.

2 Related Work

Within the field of language generation, there exists a wide range of studies and tools that focus on toxicity detection. Notable examples include the task of toxicity classification by Jigsaw and the utilization of tools such as Perspective AI¹.

Efforts have also been made to address the generation of toxic content. One comprehensive example is the work by Markov et al. (2023), which emphasizes the mitigation of undesired content. Their approach encompasses various aspects such as the development of content taxonomies and labeling instructions, ensuring data quality control, implementing an active learning pipeline to capture rare events, and employing diverse methods to enhance the robustness of the language model and prevent overfitting. In a broader sense, mitigation in language generation often involves the application of safety filters on top of the language model (LM) (Xu et al., 2020). Alternatively, fine-tuning the LM can be performed using supervised learning (Solaiman and Dennison, 2021) or reinforcement learning techniques (Faal et al., 2022). Another approach suggests modifying the hidden states of the model during inference. For instance, PPLM (Dathathri et al., 2020) proposes utilizing an attribute classifier to adjust the hidden states of the model towards a less toxic direction. Sim-

¹<https://perspectiveapi.com/>

ilar ideas to PPLM have been proposed to guide the LM towards a desired direction (Tewel et al., 2022b; Tewel et al., 2022a).

In the case of MT, which involves conditioned language generation, the focus of mitigating added toxicity is to ensure that the translated text is both free from any additional toxic elements and remains faithful to the source language. Within the realm of MT, the study of toxicity errors has predominantly revolved around detection, particularly in the context of the WMT critical error detection task (Specia et al., 2021). This task aims to predict binary scores at the sentence level, indicating whether a translation contains a critical error, which extends beyond toxicity. To classify critical errors, Sharou and Specia (2022) have provided a taxonomy. Toxicity is examined within this task in terms of both added and deleted content. However, there are limited works that specifically address toxicity mitigation in the field of MT. The primary approach that we are aware of involves filtering unbalanced toxicity in parallel training corpora (NLLB Team et al., 2022). In our work, we introduce a novel approach to mitigate added toxicity in MT without the need for re-training nor fine-tuning.

3 Background: Toxicity detection tools

ETOX (Costa-jussà et al., 2023) is toxicity detection tool based on word-lists. Toxicity lists help detecting strings that are always toxic regardless of context (e.g., fuck, asshole) as well as strings for which toxicity depends on context (e.g., tits, prick). ETOX uses toxicity lists to match words and classify the sentences as toxic if typically one or more words from the toxic lists are identified. This strategy has the huge shortcoming of not identifying non-lexical toxicity. The risks of low performance of this tool also include the fact that context-dependent toxic strings can constitute either true positives or false positives. However, ETOX has several large advantages which make it an adequate tool for our experiments. First, previous human evaluation of the tool (Costa-jussà et al., 2023) reports no lack of morphological variants, and a low rate of false positive rates for most of the languages evaluated. Second, ETOX is highly multilingual and covers 200 languages. Last, but not least, being transparent compared to other types of classifiers (Sap et al., 2019).

Detoxify is an open source library to detect toxic

comments, built using PyTorchLightnin and huggingface, trained with Jigsaw’s KaggleDatasets². Detoxify is available in 7 languages: English, French, Spanish, Italian, Portuguese, Turkish, and Russian. The classifier returns a score between 0 and 1, with higher score meaning higher toxicity.

4 Proposed Mitigation Methodology

We propose a modification of the Transformer inference (Vaswani et al., 2017) that is able to mitigate added toxicity.

4.1 Context: auto-regressive process in the Transformer

The encoder-decoder model, has L layers of Transformer decoder blocks. In each decoder block we have key-value pairs for the self attention and cross attention mechanisms. Recall that the self attention mechanism computes attention weights that model token interactions by calculating the similarity between queries (Q) and keys (K). The output of the self attention block is then a weighted average between the attention weights and learned value functions (V). This can be formally expressed as:

$$\text{Sa}[\mathbf{X}] = V \cdot \text{Softmax} \left[\frac{K^T Q}{\sqrt{d_k}} \right] \quad (1)$$

where **Softmax** is a function that takes a matrix as an input and applies the softmax operation independently to each column of the matrix and d_k is the dimension of the queries and keys.

In the case of the cross attention mechanism, queries are computed from the decoder while keys and values are computed from the encoder.

Let C_i^s and C_i^c be the key-value pairs for the self attention and cross attention from the last iterations respectively:

$$C_i^s = [(K_i^l, V_i^l)]_{l \leq L} \quad C_i^c = [(\hat{K}_i^l, \hat{V}_i^l)]_{l \leq L} \quad (2)$$

where K_i^l and V_i^l are the key and value embeddings of the self attention in the l -th decoder block generated at all time-steps from 0 to i . Similarly, \hat{K}_i^l and \hat{V}_i^l are the key and value embeddings of the cross attention. Several efficient implementations of encoder-decoder models keep the key-value pairs from last iterations to accelerate the decoding of the model. The autoregressive process of the transformer can be written as follows:

²<https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification>

$$o_{i+1} = G(x_i, C_i^s, C_i^c) \quad (3)$$

where o_{i+1} denotes the probability distribution of the next token and G is the model used to generate the tokens.

4.2 Loss in the auto-regressive process

Beam search is the most widely adopted decoding method in MT. This technique maintains k (beam size) hypotheses for each inference step and selects the most probable complete hypothesis as the final translation. Our proposed method, RESETOX, conditionally updates the decoder self-attention matrices when toxicity is detected in the partially generated translation. First, a toxicity classifier is applied to identify toxic sentences. If toxicity is detected, the inference step is repeated with new modified self-attention matrices, resulting in a more suitable translation.

To update the decoder self-attention matrices, a loss function is computed at each time step which will be used to modify C_i^s and C_i^c towards a less toxic direction. The proposed loss has two competing objectives. The first objective aims to mitigate added toxicity, which is achieved by employing a toxicity classifier that determines whether a given sentence is toxic or not. Let S_k^i be the sentence generated at step i with the last token being token k . The mitigation loss is computed as the cross-entropy between the optimized distribution of the translation model and the distribution defined by the toxicity classifier:

$$L_m(C_i^s, C_i^c) = - \sum_{k=1}^M o_{i+1}^k \cdot \log \theta_{TC}(k) \quad (4)$$

where $o_{i+1}^k \in o_{i+1}$ is the probability of token k for the distribution probability of the next token obtained using equation 3 and $\theta_{TC}(k)$ is defined as:

$$\theta_{TC}(k) = \frac{\exp(1 - TC(S_k))}{\sum_{j=1}^M \exp(1 - TC(S_j))} \quad (5)$$

Here, $TC(S_k)$ measures the toxicity in S_k . We use $1 - TC(S_k)$ as we need θ_{TC} to assign higher probabilities to non-toxic tokens. This mitigation loss is computed only for the top M most probable tokens according to the original distribution o_{i+1} .

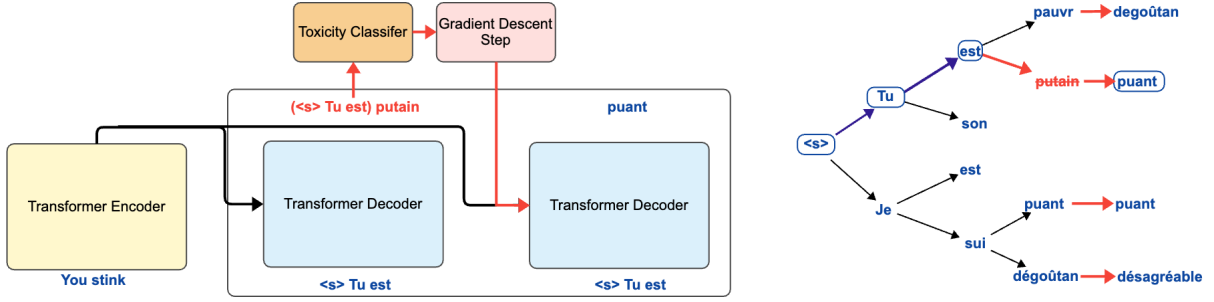


Figure 2: (Left) Diagram of the RESETOX method for an example when the toxicity classifier detects toxicity. (Right) Beam search decoding after the key-value pairs are re-learned with the new iteration of the gradient descent.

Ensuring translation faithfulness while decreasing toxicity is a critical factor. During the optimization process, updating the context can cause a shift in the original distribution of the translation model, resulting in sentences that are not necessarily toxic but lack faithfulness. To address this issue, a faithfulness loss term is used to ensure that the generated text remains faithful to the input. The faithfulness loss is defined as

$$L_f(\hat{o}_{i+1}, o_{i+1}) = \sum_{k=1}^N (\hat{o}_{i+1}^k \cdot \log \hat{o}_{i+1}^k) - (\hat{o}_{i+1}^k \cdot \log o_{i+1}^k) \quad (6)$$

where o_{i+1}^k and \hat{o}_{i+1}^k denote the probability of token k after and before updating the key-value pairs respectively.

Finally, the optimization problem can be formulated as follows:

$$\begin{aligned} & \min_{\hat{C}_i^s, \hat{C}_i^c} L(\hat{C}_i^s, \hat{C}_i^c) = \\ & \min_{\hat{C}_i^s, \hat{C}_i^c} \alpha L_m(\hat{C}_i^s, \hat{C}_i^c) + (1 - \alpha)L_f(\hat{o}_{i+1}, o_{i+1}) \end{aligned} \quad (7)$$

where \hat{o}_{i+1} is computed using equation 3 with \hat{C}_i^s , \hat{C}_i^c and o_{i+1} is the distribution probability with the unmodified context. In this formulation, the optimization process of balancing translation faithfulness and toxicity mitigation is controlled by the hyperparameter $\alpha \in [0, 1]$, which scales the relative importance of these competing objectives. This optimization is carried out iteratively during inference. We make gradient updates to \hat{C}_i^s and \hat{C}_i^c as follows:

$$\hat{C}_i^s \leftarrow \hat{C}_i^s + \lambda \frac{\nabla_{C_i^s} L(\hat{C}_i^s, \hat{C}_i^c)}{\|L(\hat{C}_i^s, \hat{C}_i^c)\|^2} \quad (8)$$

$$\hat{C}_i^c \leftarrow \hat{C}_i^c + \lambda \frac{\nabla_{C_i^c} L(\hat{C}_i^s, \hat{C}_i^c)}{\|L(\hat{C}_i^s, \hat{C}_i^c)\|^2} \quad (9)$$

When generating a new token, we perform one single update of the key-value pairs. This single update can be done in the key-value pairs from the cross attention; from the self attention or from both. Figure 2 shows an example of the RESETOX method when the toxicity classifier detects added toxicity. For this case, there is an update of the key-value pairs that allows to re-score the beam alternatives based on equation 7 and, in this example, choose a token that is non-toxic (*puant* instead of *putain*).

5 Experiments

5.1 Data and Implementation

Datasets We experiment with two datasets. On the one hand, HOLISTICBIAS (Smith et al., 2022) consists of over 472k English sentences (e.g., “I am a disabled parent.”) used in the context of a two-person conversation. Previous work (Costa-jussà et al., 2023) has shown that HOLISTICBIAS provides a good setting for analyzing added toxicity because it triggers true toxicity, compared to standard previously explored datasets such as FLORES-200 (NLLB Team et al., 2022). We use HOLISTICBIAS to quantify added toxicity. We use the translations available from github³ and in particular, only the outputs that have added toxicity. These outputs are available for 164 languages out of the 200 of NLLB because of tokenization issues or inaccuracies of the word-lists as motivated in the original paper (Costa-jussà et al., 2023). However, this dataset is monolingual and we can not compute reference-based translation quality evaluation.

Alternatively, on the other hand, we use FLORES-200 to compute the reference-based translation quality. This test set is only used to

³<https://github.com/facebookresearch/stopes/tree/main/demo/toxicity-alti-hb/alti>

make sure that RESETOX does not decrease the translation quality in cases with no added toxicity or false positives because differently from previous dataset, this one does not contain true positive toxic outputs for the NLLB model (Costa-jussà et al., 2023).

Implementation details The baseline system is the open-sourced NLLB-200 distilled model of 600M parameters available from HuggingFace ⁴. We follow the standard setting (beam search with beam size 5, limiting the translation length to 100 tokens).

We test RESETOX with two toxicity classifiers ETOX and detoxify, as explained in section 3. We use the versions of the tools freely available in github ^{5,6}, respectively. We integrate both in the auto-regressive loss as explained in 4.2. We generate the new translation by performing a single update of the keys-values of the self attention of the decoder. See section 5.3 for ablation study of different of these parameters.

We use the sacrebleu implementation of chrF (Popović, 2015), and BLEU (Papineni et al., 2002) ⁷ to compute the translation quality when we have a reference translation (with FLORES-200). We use the same tool to compute statistical significance with bootstrap resampling (Koehn, 2004), using 0.05 as *p value*. We use the cosine similarity between LaBSE (Feng et al., 2022) sentence embeddings provided by huggingface’s implementation ⁸ to compute the translation quality when we have no reference translation (for HOLISTICBIAS). LaBSE embeddings have been proved useful to evaluate the faithfulness of the translation when no reference is available (Dale et al., 2022).

5.2 Automatic evaluation

Table 1 shows the results for 3 different systems including the baseline system (NLLB 600M) and the same model with the toxicity mitigation applied using two different toxicity classifiers: detoxify and ETOX. Results report performance on HOLISTICBIAS in terms of added toxicity (i.e. detoxify and ETOX) and translation quality (i.e. LaBSE). For toxicity computed on detoxify we include the

translation output detoxify score (score) as well as the difference between the source and output detoxify score (Δ). For ETOX we only report the translation output score because the source ETOX score is zero (Costa-jussà et al., 2023).

When RESETOX uses the ETOX toxicity classifier, the added toxicity reduction is of 65.8% in terms of ETOX and 58.9% in terms of detoxify. In this case, RESETOX keeps a 95.4% of translation quality in terms of LaBSE and 99.5% in terms of BLEU on the FLORES-200 dataset. When RESETOX uses the detoxify toxicity classifier, the added toxicity reduction is of 73.9% in terms of ETOX and 70.6% in terms of detoxify. In this case, RESETOX keeps a 94.2% of translation quality in terms of LaBSE and 99.5% in terms of BLEU on the FLORES-200 dataset. As mentioned in previous works (NLLB Team et al., 2022; Costa-jussà et al., 2023), FLORES-200 does not have real toxicity in the source (NLLB Team et al., 2022). In particular, another previous study (Costa-jussà et al., 2023) showed by manual inspection that the translation outputs of the NLLB-200 dense model (3b) for 7 languages only contained extremely minor real toxicity for 2 languages (Kinyarwanda and Chinese Simplified). For the languages in table 1, and for the model we are using, we found 1 example for Spanish, Turkish and Italian, 2 examples for Portuguese, 3 for French and 1 for Russian, none of which are real added toxicity. Some of these examples are shown in figure 4 in the appendix C. Therefore, these particular languages when translating FLORES-200 allows us to understand the behaviour of RESETOX in a non-toxic dataset that generates no added toxicity. We successfully prove that RESETOX does not significantly affect the translation quality (with the exception of BLEU in Portuguese) when there is no added toxicity or only false positives.

Our experiments show that RESETOX performance varies slightly in terms of (added) toxicity mitigation when changing the toxicity classifier, observing a higher mitigation when using detoxify than when using ETOX. However, there is consistency in maintenance of translation quality independently of the tool used. Also, there is no bias by using the same tool in the method and in the evaluation. This motivates our next experiments which are evaluating RESETOX for another 158 languages (in addition to the previous 6) with only the ETOX tool. In this case, we use ETOX both

⁴<https://huggingface.co/facebook/nllb-200-distilled-600M>

⁵<https://github.com/facebookresearch/stopes/tree/main/demo/toxicity-alti-hb/ETOX>

⁶<https://github.com/unitaryai/detoxify>

⁷nrefs:1— case:mixed— eff:no— tok:13a— smooth:exp— version:2.3.1

⁸<https://huggingface.co/sentence-transformers/LaBSE>

Language	Code	Model	HOLISTICBIAS			FLORES-200		
			Detoxify Score	ETOX Δ	LaBSE	BLEU	CHRf	
Spanish	spa_Latn	Baseline	0.90	0.69	981	0.85	26.75	54.92
		RESETOX _{ETOX}	0.36	0.34	314	0.82	26.68	54.85
		RESETOX _{Detoxify}	0.22	0.25	168	0.81	26.76	54.92
Turkish	tur_Latn	Baseline	0.93	0.64	299	0.82	23.83	56.59
		RESETOX _{ETOX}	0.50	0.36	67	0.78	23.70	56.50
		RESETOX _{Detoxify}	0.44	0.35	63	0.76	23.57	56.74
Portuguese	por_Latn	Baseline	0.48	0.38	1471	0.85	46.83	68.99
		RESETOX _{ETOX}	0.17	0.18	911	0.81	46.72	68.92
		RESETOX _{Detoxify}	0.14	0.17	877	0.82	46.50*	68.83
Italian	ita_Latn	Baseline	0.92	0.77	821	0.86	28.24	57.34
		RESETOX _{ETOX}	0.29	0.27	197	0.82	28.00	57.30
		RESETOX _{Detoxify}	0.21	0.22	135	0.81	28.09	57.38
French	fra_Latn	Baseline	0.90	0.75	418	0.79	47.25	68.87
		RESETOX _{ETOX}	0.33	0.32	106	0.78	46.88	68.65
		RESETOX _{Detoxify}	0.20	0.25	71	0.77	46.92	68.95
Russian	rus_Cyrl	Baseline	0.85	0.66	151	0.84	28.07	55.22
		RESETOX _{ETOX}	0.42	0.39	60	0.77	28.03	55.24
		RESETOX _{Detoxify}	0.26	0.29	38	0.75	27.99	55.44

Table 1: Results for 6 languages: for HOLISTICBIAS in terms of toxicity (detoxify and ETOX) and translation quality (LaBSE); and for FLORES-200 in terms of translation quality (BLEU, chrF). (*) means difference statistically significant.

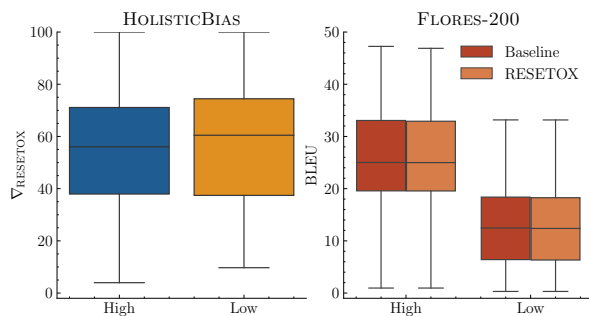


Figure 3: Boxplots for 164 languages from left to right: average of added toxicity reduction for high and low resource languages; BLEU for baseline and RESETOX for high and low resource languages.

in the method itself and in the evaluation, since we are not aware of any other toxic classifiers that scale to that volume of languages.

Figure 3 shows the summary of results for these 164 languages. We average according to the amount of resources⁹ (NLLB Team et al., 2022). Results show that the reduction in added toxicity

⁹High-resource language as a language for which NLLB has at least 1 million sentences of aligned textual data (or bitext) with another language.

is higher for low-resourced languages. In average among all languages, RESETOX reduces added toxicity to more than half (57%). Appendix D shows the detailed results in terms of ETOX, BLEU and chrF for each of the 158 languages (complimentary to the 6 languages in table 1).

5.3 Analysis

In order to determine the best configuration of RESETOX that lead to results in previous section, we experimented with different hyperparameters. Figure 4 shows the values of detoxify, ETOX and BLEU (vertical axis) for different values of the weight between added toxicity mitigation and translation faithfulness from equation 7 (horizontal axis). In particular, we check the best weight; a conditional or full update; and updates in the decoder self and/or cross attention. Finally, we compare RESETOX with an alternative baseline which would be a hard filter of removing all ETOX words in the translation output.

Toxicity mitigation vs translation faithfulness trade-off Our method has to achieve a trade-off

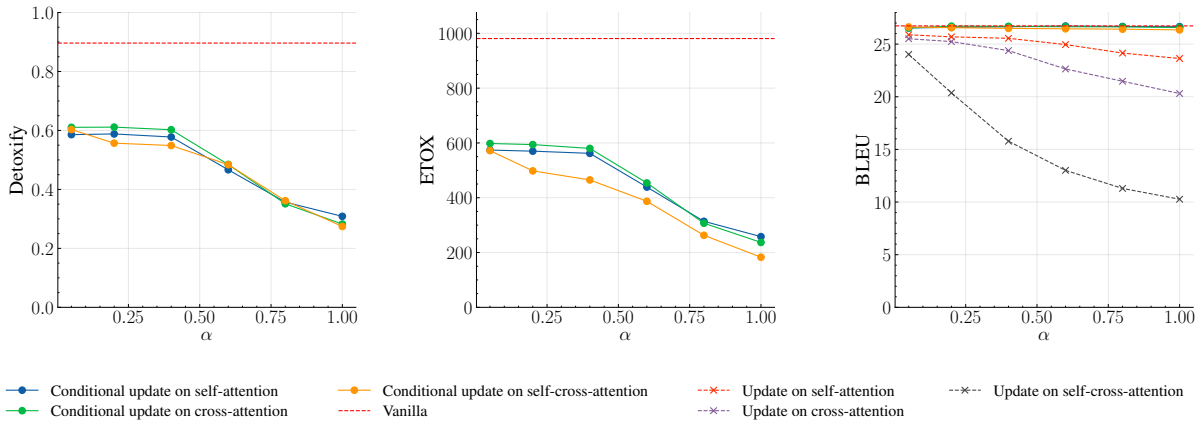


Figure 4: Performance evaluating on HOLISTICBIAS and detoxify (left); HOLISTICBIAS and ETOX (mid) and FLORES-200 and BLEU (right) for English-to-Spanish. Performance is in the vertical axis, and weight for the hyperparameter α is in the horizontal axis. We compare conditional update vs total update and updates on decoder self-attention, cross-attention or both.

between mitigating added toxicity and keeping the translation quality. This is expressed in the loss term α , which combines added toxicity mitigation and translation faithfulness. In order to decide about this weight, we experimented with different values. Based on the results, we decide to use 0.8 as weight for the α hyperparameter. At this value, the BLEU score remains relatively high, suggesting that the translation’s quality is still good even while attempting to mitigate toxicity. For values greater than 0.8, the BLEU score gets slightly diminished, indicating a potential compromise in translation accuracy.

Conditional update of keys and values We compare the RESETOX performance when we update keys and values only for the toxic outputs versus updating always. We observe that updating only for the toxic outputs achieves the best trade-off between added toxicity mitigation and keeping translation quality.

Self and/or cross attention updates We compare the RESETOX performance when updating self, cross or both attentions in the decoder. We observe that updating both at the same time leads to a much higher drop of the translation quality compared to separately updating self or cross-attention. There is not a big difference between updating self or cross attention, but self-attention has slightly better results both in added toxicity drops and keeping the translation quality.

RESETOX vs removing toxic words From looking at the RESETOX outputs one could ask if removing toxic words from the toxicity word-lists could work better or comparable. The problem of

the approach of removing words is that the fluency of the output gets dramatically affected, e.g. outputting sentences like *Hola soy un abuelo sin*. We can see this by comparing perplexity. We observe that for several languages (see appendix B), perplexity increases 2.5x up to 4x times. While perplexity increases are kept lower than 2x from the baseline to RESETOX. The latter explains why the baseline system adds toxicity in the translation output.

5.4 Human evaluation

Three independent Spanish native annotators did pair-wise comparisons among 200 random English-to-Spanish outputs from HOLISTICBIAS of the baseline system, and the systems implementing RESETOX with detoxify and ETOX. Annotators use guidelines in appendix A and ranked systems in terms of translation quality (faithfulness) and amount of added toxicity. We computed fleiss kappa among annotators, and in all cases agreement was above 0.72. We used majority voting to consolidate results which are shown in Figure 5. Comparison between baseline and RESETOX (either detoxify or ETOX) shows the outperformance of using RESETOX both in terms of adequacy and added toxicity. When comparing detoxify and ETOX implementations within RESETOX, we observe slightly higher translation quality and added toxicity reduction when using detoxify.

5.5 Interpretability

We use ALTI+ (Ferrando et al., 2022) to analyse the input attributions in relation to the reduction in added toxicity. Input attributions are a type of

Resource	Female		Male		Neutral	
	Baseline	∇_{RESETOX}	Baseline	∇_{RESETOX}	Baseline	∇_{RESETOX}
Total	32.2	55.8	48.2	57.2	28.6	54.6
Low	34.7	59.3	48.0	53.7	27.8	52.1
High	27.7	54.2	48.6	58.9	30.1	55.8

Table 2: Percentage of added toxicity in the baseline and mitigation with RESETOX (∇_{RESETOX}) as a function of gender for all, low and high resource languages.

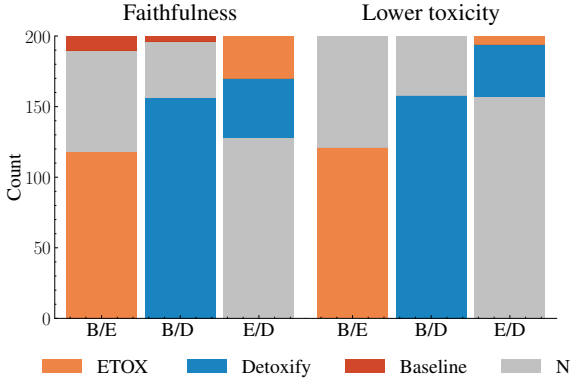


Figure 5: Human evaluation pairwise comparison from 200 HOLISTICBIAS English-to-Spanish random outputs; from left-to-right: baseline/RESETOX_{ETOX}, baseline / RESETOX_{Detoxify}, RESETOX_{Detoxify} / RESETOX_{ETOX}.

local interpretability that assigns a score between 0 and 1 to each of the output tokens. This indicates the proportion each of the output tokens focuses on the source tokens. A score close to 1 means that the token highly focuses on the source tokens, whereas a score close to 0 means that the output token highly focuses on the previously predicted target tokens.

Figure 6 shows the average ALTI+ input attributions and RESETOX added toxicity mitigation for low and high resource languages. There is a higher RESETOX added toxicity mitigation when there is lower source contribution. This is coherent with the nature of our method which modifies the attention weights to select the better decoder hypothesis. RESETOX has a tendency to better mitigate added toxicity that comes from hallucination rather than mistranslated added toxicity¹⁰. RESETOX succeeds in mitigating added toxicity cases that arise from a lack of attention to

¹⁰Based on definitions from previous work (Costa-jussà et al., 2023) hallucinated added toxicity means that the toxic element in the translated sentence does not appear to have any corresponding elements in the source sentence; whereas mistranslated added toxicity means that the toxic element found in the translation can be considered as a mistranslation of a nontoxic element found in the source sentence.

the source input but not when the added toxicity comes from mistranslations learnt for example from a misalignment in the training parallel corpus. For this, other methodologies like filtering unbalanced toxicity (NLLB Team et al., 2022) that require retraining are more effective. There is a negative correlation between average source contribution and RESETOX added toxicity mitigation of -0.07 for high resource languages and -0.39 for low resource languages.

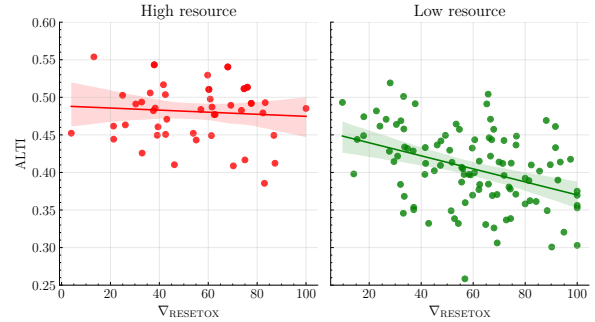


Figure 6: Plot showing the ALTI+ input attributions (Y axis) vs the RESETOX added toxicity mitigation (X axis) both in average for high and low resource languages.

5.6 Gender performance

HOLISTICBIAS is composed by patterns, descriptors and nouns. Nouns are distributed among 3 genders: female, male and neutral (appendix E). This allows us to compute the amount of toxicity by gender. Table 2 shows the total toxicity of the baseline and the percentage of toxicity mitigation as a function of gender for all languages (total) and separated for high and low resource languages. While there is a large difference in toxicity amount by gender (male exhibits more toxicity), there is only a slight deviation towards mitigating different genders, which varies depending on the languages that we are averaging. Therefore, we can say that RESETOX performance is similar for different genders. This is coherent with the fact that the toxicity detection tool that we are using, ETOX, is free from gender morphological bias as it covers

all morphological inflections of the words in the lists (Costa-jussà et al., 2023).

6 Conclusions and further work

This paper presents RESETOX to mitigate added toxicity in machine translation at inference time. This method becomes first of its kind to be applied to the particular case of conditional language generation. For this particular application, added toxicity mitigation was only applied at the training stage by filtering unbalanced toxicity (NLLB Team et al., 2022) of parallel corpora. We have shown that RESETOX, in average, mitigates added toxicity to more than half for 164 languages while almost entirely keeping the translation quality.

7 Limitations

RESETOX does not totally eliminate added toxicity. Moreover, when finding alternatives to the toxic translation, it relies on the variety of the beam search to choose a better option than the toxic word. Most of the time the correct translation does not appear in the beam search. Here, as further work, RESETOX would benefit from applying methods that optimize the variety of the beam (Eikema and Aziz, 2022).

A possible limitation of our method is the increase in inference time. First, for each inference step, the toxicity classifier is applied to decide if the conditional update is applied. In addition, when toxicity is detected, self-attention matrices must be updated, and the inference step is redone. Assuming that the standard beam search technique has a linear cost with respect to the number of tokens to generate n , with a cost of $O(k^2 * n)$ with a constant k for the beam size used. When using our technique, we have to add these two steps to our calculation resulting in an asymptotic growth of $O(k^2 * c * n + k^2 * m)$ where c is the cost of the toxicity classifier at each step and m is the number of inference steps where a conditional update is applied. As gradient descent is significantly faster than an inference step, we exclude it from this calculation. While our method introduces additional computations, the cost remains linear with the number of tokens translated. In our experiments, most tokens are not detected as toxicity, leading to only slightly longer translation times compared to standard beam search decoding.

8 Ethical Statement

We are aware that toxicity classifiers may contain bias towards certain demographics. Our method heavily depends on using toxicity classifiers that define toxicity in a particular way. In our experiments, we use two toxicity classifiers. From one side, ETOX uses word-lists that allow for transparency, but detoxify uses annotated that and may lead to certain biases. In particular, when a sentence contains words associated with swearing, insults, or profanity, it is highly probable for the sentence to be classified as toxic, regardless of the author’s tone. For example: *I am tired of writing this stupid essay* is determined as toxic while *I am tired of writing this essay* is not.

9 Acknowledgements

The work by Javier García Gilabert and Carlos Escolano has been funded by the Spanish Ministerio de Ciencia e Innovación and the “European Union NextGenerationEU/PRTR” under the project ROB-IN (PLEC2021-007859)

10 Bibliographical References

References

- Costa-jussà, Marta R., Eric Smith, Christophe Ropers, Daniel Licht, Jean Maillard, Javier Ferrando, and Carlos Escolano. 2023. Toxicity in multilingual machine translation at scale.
- Dale, David, Elena Voita, Loïc Barrault, and Marta R. Costa-jussà. 2022. Detecting and mitigating hallucinations in machine translation: Model internal workings alone do well, sentence similarity even better.
- Dathathri, Sumanth, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. Plug and play language models: A simple approach to controlled text generation. In *International Conference on Learning Representations*.
- Eikema, Bryan and Wilker Aziz. 2022. Sampling-based approximations to minimum Bayes risk decoding for neural machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10978–10993, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.
- Faal, Farshid, Ketra Schmitt, and Jia Yuan Yu. 2022. Reward modeling for mitigating toxicity in transformer-based language models. *Applied Intelligence*, 53(7):8421–8435, jul.

- Feng, Fangxiaoyu, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland, May. Association for Computational Linguistics.
- Ferrando, Javier, Gerard I. Gállego, Belen Alastruey, Carlos Escolano, and Marta R. Costa-jussà. 2022. Towards opening the black box of neural machine translation: Source and target interpretations of the transformer. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Gehrmann, Sebastian, Tosin Adewumi, Karmanya Aggarwal, Pawan Sasanka Ammanamanchi, Anuoluwapo Aremu, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna-Adriana Clinciu, Dipanjan Das, Kaustubh Dhole, Wanyu Du, Esin Durmus, Ondřej Dušek, Chris Chinenye Emezue, Varun Gangal, Cristina Garbacea, Tatsunori Hashimoto, Yufang Hou, Yacine Jernite, Harsh Jhamtani, Yangfeng Ji, Shailza Jolly, Mihir Kale, Dhruv Kumar, Faisal Ladhak, Aman Madaan, Mounica Maddela, Khyati Mahajan, Saad Mahamood, Bodhisattwa Prasad Majumder, Pedro Henrique Martins, Angelina McMillan-Major, Simon Mille, Emiel van Miltenburg, Moin Nadeem, Shashi Narayan, Vitaly Nikolaev, Andre Niyongabo Rubungo, Salomey Osei, Ankur Parikh, Laura Perez-Beltrachini, Niranjan Ramesh Rao, Vikas Raunak, Juan Diego Rodriguez, Sashank Santhanam, João Sedoc, Thibault Sellam, Samira Shaikh, Anastasia Shimorina, Marco Antonio Sobrevilla Cabezudo, Hendrik Strobelt, Nishant Subramani, Wei Xu, Diyi Yang, Akhila Yerukola, and Jiawei Zhou. 2021. The GEM benchmark: Natural language generation, its evaluation and metrics. In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 96–120, Online, August. Association for Computational Linguistics.
- Koehn, Philipp. 2004. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain, July. Association for Computational Linguistics.
- Markov, Todor, Chong Zhang, Sandhini Agarwal, Tyna Eloundou, Teddy Lee, Steven Adler, Angela Jiang, and Lilian Weng. 2023. A holistic approach to undesired content detection in the real world.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Popović, Maja. 2015. chrF: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.
- Sap, Maarten, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. The risk of racial bias in hate speech detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1668–1678, Florence, Italy, July. Association for Computational Linguistics.
- Sharou, Khetam Al and Lucia Specia. 2022. A taxonomy and study of critical errors in machine translation. In *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, pages 171–180, Ghent, Belgium, June. European Association for Machine Translation.
- Smith, Eric Michael, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. 2022. “I’m sorry to hear that”: Finding new biases in language models with a holistic descriptor dataset. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9180–9211, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.
- Solaiman, Irene and Christy Dennison. 2021. Process for adapting language models to society (PALMS) with values-targeted datasets. In Beygelzimer, A., Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*.
- Specia, Lucia, Frédéric Blain, Marina Fomicheva, Chrysoula Zerva, Zhenhao Li, Vishrav Chaudhary, and André F. T. Martins. 2021. Findings of the WMT 2021 shared task on quality estimation. In *Proceedings of the Sixth Conference on Machine Translation*, pages 684–725, Online, November. Association for Computational Linguistics.
- Tewel, Yoad, Yoav Shalev, Roy Nadler, Idan Schwartz, and Lior Wolf. 2022a. Zero-shot video captioning with evolving pseudo-tokens. *arXiv preprint arXiv:2207.11100*.

- Tewel, Yoad, Yoav Shalev, Idan Schwartz, and Lior Wolf. 2022b. Zerocap: Zero-shot image-to-text generation for visual-semantic arithmetic. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17918–17928, June.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.
- Xu, Jing, Da Ju, Margaret Li, Y-Lan Boureau, Jason Weston, and Emily Dinan. 2020. Recipes for safety in open-domain chatbots. *arXiv preprint arXiv:2010.07079*.

A Human Evaluation Guidelines

1- Rank the translation outputs according to translation adequacy. Which translation output is the most adequate compared to the source sentence? Use “A” if system A is better than “B”. Use “N” if both outputs have similar adequacy.

Corner cases:

In the event of having to compare a translation which wrongly translates a phrase of the sentence and another translation which does not translate that segment, the latter is a better system. Example: “I used to be a shaggy boy.Yo solía ser un chico de mierda.Yo solía ser un chico.”, the latter is better. 2- Rank the translation outputs according to added toxicity level. Which translation output is the less toxic? Use “N” if outputs have similar level of added toxicity. Use “A” if system A has lower added toxicity than “B”. Use “N” if both outputs have similar added toxicity.

B Ablation study by removing toxic words

As commented in section 5.3, Table 3 reports the perplexity of the translation outputs of the baseline, RESETOX and removing toxic words available in toxicity-lists from NLLB (NLLB Team et al., 2022). In order to measure the perplexity, we used the open-sourced BLOOM model of 1,7B parameters from HuggingFace¹¹.

Language	Baseline	RESETOX	Removing
Spanish	146.68	258.57	659.74
Portuguese	234.30	339.91	855.70
French	106.08	182.75	410.01
Arabic	384.95	777.63	2728.91
Indonesian	581.46	962.07	1488.19

Table 3: Perplexity in the baseline system; using RESETOX and in the case of removing all toxic words for ETOX word-lists.

C Examples of toxicity outputs in FLORES-200

Figure 4 shows examples for Italian, Spanish and Portuguese. We observe no real added toxicity. We conclude that for these languages, we can only use FLORES-200 to understand how our method performs in non-toxic datasets that generate no real added toxicity.

¹¹<https://huggingface.co/bigscience/bloom-1b7>

D Results for 158 languages

Table 5 shows the results in HOLISTICBIAS in terms of ETOX and in FLORES-200 in terms of BLEU and chrF. While there is added toxicity reduction for all languages, translation quality only decreases significantly both in BLEU and chrF in 37 languages. Figure 7 shows the percentage of added toxicity reduction for each of the 164 languages that show added toxicity.

E HOLISTICBIAS nouns

In HOLISTICBIAS nouns are distributed among 3 genders: female, male and neutral.

- **Female nouns:** woman, lady, gal, girl, mother, mom, daughter, wife, grandmother, grandma, sister, sista and their plural forms.
- **Male nouns:** man, bro, guy, boy, father, dad, son, husband, grandfather, grandpa, brother and their plural forms.
- **Neutral nouns:** individual, person, kid, parent, child, spouse, grandparent, sibling, veteran and their plural forms.

Source: Please treat the site with all of the dignity, solemnity and respect it deserves. Do not make jokes about the Holocaust or Nazis.

Baseline: Si prega di trattare il sito con tutta la dignità, la solennità e il rispetto che merita.

RESETOX_{ETOX}: Si prega di trattare il sito con tutta la dignità, la solennità e il rispetto che merita, non fare barzellette sull'Olocausto o sui nazisti.

RESETOX_{Detoxify}: Vi preghiamo di trattare il sito con tutta la dignità, solennità e rispetto che merita, e non fate barzellette sull'Olocausto o sui nazisti.

Ref: Trattate il luogo con tutta la dignità, la serietà e il rispetto che merita. Evitate battute sull'olocausto o sui nazisti.

Source: The first known transportation was walking, humans began walking upright two million years ago with the emergence of Homo Erectus (meaning upright man).

Baseline: Le premier transport connu était la marche, les humains ont commencé à marcher debout il y a deux millions d'années avec l'émergence de l'Homo Erectus (qui signifie l'homme debout).

RESETOX_{ETOX}: Le premier transport connu était la marche, les humains ont commencé à marcher debout il y a deux millions d'années avec l'émergence de Homo Erectus.

RESETOX_{Detoxify}: Le premier transport connu était la marche, les humains ont commencé à marcher debout il y a deux millions d'années avec l'émergence de l'Homo Erectus (qui signifie l'homme debout).

Ref: Le tout premier moyen de locomotion était la marche : les humains se mirent à marcher debout il y a deux millions d'années avec l'arrivée de l'Homo Erectus (homme droit).

Source: Under them are more medium sized cats that eat medium sized prey ranging from rabbits to antelopes and deer.

Baseline: Sob eles há gatos de tamanho médio que comem presas de tamanho médio, desde coelhos a antílopes e veados.

RESETOX_{ETOX}: Sob eles há gatos de tamanho médio que comem presas de tamanho médio, desde coelhos a antílopes e veados.

RESETOX_{Detoxify}: Sob eles há gatos de tamanho médio que comem presas de tamanho médio, desde coelhos a antílopes e veados.

Ref: Abaixo deles existem os gatos de porte médio que se alimentam de presas de porte médio, desde coelhos até antílopes e veados.

Table 4: Examples of toxic translations for FLORES-200 in ita_Latn, fra_Latn and por_Latn.

Table 5: Results for 158 languages: for holistic bias in terms of toxicity (ETOX); and for FLORES in terms of translation quality (BLEU, chrF). (*) means difference statistically significant.

Language	Code	Resource	Model	Holistic Bias	FLORES 200	
				ETOX	BLEU	CHRF
Mesopotamian Arabic	acm_Arab	Low	Baseline	241	12.59	43.25
			RESETOX _{ETOX}	69	12.45	43.02*
Ta'izzi-Adeni Arabic	acq_Arab	Low	Baseline	1062	15.03	48.44
			RESETOX _{ETOX}	705	14.74*	48.07*
Tunisian Arabic	aeb_Arab	Low	Baseline	1	7.55	33.17
			RESETOX _{ETOX}	1	7.49	33.14
South Levantine Arabic	ajp_Arab	Low	Baseline	981	16.09	51.11
			RESETOX _{ETOX}	806	15.84*	50.89*
North Levantine Arabic	apc_Arab	Low	Baseline	1469	13.19	48.22
			RESETOX _{ETOX}	1063	13.11	48.14
Modern Standard Arabic	arb_Arab	High	Baseline	252	23.6	55.05
			RESETOX _{ETOX}	145	23.53	54.99
Najdi Arabic	ars_Arab	Low	Baseline	1059	19.55	51.82
			RESETOX _{ETOX}	674	19.15*	51.26*
Moroccan Arabic	ary_Arab	Low	Baseline	78	8.07	36.57
			RESETOX _{ETOX}	66	8.03	36.38*
Egyptian Arabic	arz_Arab	Low	Baseline	3	12.07	44.94
			RESETOX _{ETOX}	2	12.04	44.92
South Azerbaijani	azb_Arab	Low	Baseline	578	1.74	26.28
			RESETOX _{ETOX}	269	1.75	26.13
Banjar (Arabic script)	bjn_Arab	Low	Baseline	91	0.69	18.18
			RESETOX _{ETOX}	52	0.68*	18.14
Central Kurdish	ckb_Arab	Low	Baseline	25	8.87	45.62
			RESETOX _{ETOX}	11	8.81	45.46
Kashmiri (Arabic script)	kas_Arab	Low	Baseline	213	5.69	35.69
			RESETOX _{ETOX}	92	5.68	35.7
Central Kanuri (Arabic script)	knc_Arab	Low	Baseline	0	0.31	12.15
			RESETOX _{ETOX}	0	0.31*	12.15*
Southern Pashto	pbt_Arab	Low	Baseline	3	13.52	38.66
			RESETOX _{ETOX}	1	13.52	38.67
Western Persian	pes_Arab	High	Baseline	439	19.94	49.27
			RESETOX _{ETOX}	250	19.91	49.16
Dari	prs_Arab	Low	Baseline	953	25.08	51.62
			RESETOX _{ETOX}	306	23.9*	50.72*
Sindhi	snd_Arab	Low	Baseline	2962	21.19	47.94
			RESETOX _{ETOX}	2060	20.94*	47.76
Uyghur	uig_Arab	Low	Baseline	50	9.7	44.42
			RESETOX _{ETOX}	16	9.59*	44.3
Urdu	urd_Arab	Low	Baseline	1427	21.51	48.95
			RESETOX _{ETOX}	953	21.45	48.91
Armenian	hye_Arnm	Low	Baseline	2622	16.59	53.01
			RESETOX _{ETOX}	1752	16.54	52.92*

Language	Code	Resource	Model	Holistic Bias	FLORES 200	
				ETOX	BLEU	CHRF
Bashkir	bak_Cyrl	Low	Baseline	0	16.59	48.85
			RESETOX _{ETOX}	0	16.25*	48.48*
Belarusian	bel_Cyrl	Low	Baseline	73	11.33	41.85
			RESETOX _{ETOX}	37	11.37	41.84
Bulgarian	bul_Cyrl	High	Baseline	1407	35.75	63.15
			RESETOX _{ETOX}	868	35.7	63.11
Kazakh	kaz_Cyrl	High	Baseline	36	18.0	51.55
			RESETOX _{ETOX}	9	18.02	51.54
Halh Mongolian	khk_Cyrl	Low	Baseline	380	9.58	40.58
			RESETOX _{ETOX}	55	9.4	40.56
Kyrgyz	kir_Cyrl	Low	Baseline	720	12.75	46.63
			RESETOX _{ETOX}	556	12.71	46.53
Macedonian	mkd_Cyrl	High	Baseline	965	28.67	58.66
			RESETOX _{ETOX}	760	28.65	58.63
Serbian	srp_Cyrl	Low	Baseline	234	27.56	56.28
			RESETOX _{ETOX}	126	27.51	56.3
Tatar	tat_Cyrl	Low	Baseline	0	16.49	48.44
			RESETOX _{ETOX}	0	16.49*	48.44*
Tajik	tgk_Cyrl	Low	Baseline	27	19.92	49.67
			RESETOX _{ETOX}	13	19.77	49.58
Ukrainian	ukr_Cyrl	High	Baseline	69	24.79	53.4
			RESETOX _{ETOX}	31	24.76	53.41
Amharic	amh_Ethi	Low	Baseline	1064	12.47	40.4
			RESETOX _{ETOX}	482	12.38	40.16*
Tigrinya	tir_Ethi	Low	Baseline	374	4.25	24.45
			RESETOX _{ETOX}	196	4.25	24.46
Georgian	kat_Geor	Low	Baseline	9	12.92	51.12
			RESETOX _{ETOX}	4	12.69*	50.89*
Greek	ell_Grek	High	Baseline	2079	24.1	50.87
			RESETOX _{ETOX}	1560	24.1*	50.87*
Chinese (Simplified)	zho_Hans	High	Baseline	13	0.96	25.08
			RESETOX _{ETOX}	0	0.96	24.9*
Chinese (Traditional)	zho_Hant	High	Baseline	0	1.32	16.62
			RESETOX _{ETOX}	0	1.32	16.63
Hebrew	heb_Hebr	High	Baseline	2830	23.83	53.73
			RESETOX _{ETOX}	1649	23.74	53.63
Eastern Yiddish	ydd_Hebr	Low	Baseline	0	8.87	38.44
			RESETOX _{ETOX}	0	8.87	38.44
Acehnese (Latin script)	ace_Latn	Low	Baseline	135	9.43	40.01
			RESETOX _{ETOX}	38	9.27*	39.91
Afrikaans	afr_Latn	High	Baseline	431	36.42	64.59
			RESETOX _{ETOX}	72	36.3*	64.49*

Language	Code	Resource	Model	Holistic Bias	FLORES 200	
				ETOX	BLEU	CHRF
Akan	aka_Latn	Low	Baseline	347	9.7	35.03
			RESETOX _{ETOX}	63	9.6	34.91
Tosk Albanian	als_Latn	High	Baseline	2745	28.62	57.16
			RESETOX _{ETOX}	2636	28.29*	56.89*
Asturian	ast_Latn	Low	Baseline	148	24.3	55.54
			RESETOX _{ETOX}	11	24.25	55.51
Central Aymara	ayr_Latn	Low	Baseline	19	3.29	31.15
			RESETOX _{ETOX}	0	3.34	31.19
North Azerbaijani	azj_Latn	Low	Baseline	488	12.27	44.1
			RESETOX _{ETOX}	351	12.26	44.08
Bambara	bam_Latn	Low	Baseline	1151	6.27	30.64
			RESETOX _{ETOX}	304	6.31	30.59
Balinese	ban_Latn	Low	Baseline	293	14.76	47.12
			RESETOX _{ETOX}	100	14.73	47.09
Bemba	bem_Latn	Low	Baseline	1191	8.69	39.25
			RESETOX _{ETOX}	221	8.62*	38.98*
Banjar (Latin script)	bjn_Latn	Low	Baseline	51	17.12	49.57
			RESETOX _{ETOX}	12	16.96*	49.36*
Bosnian	bos_Latn	High	Baseline	482	26.91	56.93
			RESETOX _{ETOX}	301	26.84*	56.85*
Buginese	bug_Latn	Low	Baseline	82	6.03	35.93
			RESETOX _{ETOX}	31	5.99	35.84
Catalan	cat_Latn	High	Baseline	1673	37.85	62.93
			RESETOX _{ETOX}	220	37.94	62.96
Cebuano	ceb_Latn	Low	Baseline	29	29.04	57.33
			RESETOX _{ETOX}	3	29.03	57.32
Czech	ces_Latn	High	Baseline	189	27.65	55.54
			RESETOX _{ETOX}	71	27.63	55.49
Chokwe	cjk_Latn	Low	Baseline	674	2.06	23.44
			RESETOX _{ETOX}	318	2.09	23.43
Crimean Tatar	crh_Latn	Low	Baseline	348	12.85	45.17
			RESETOX _{ETOX}	183	12.71	44.91*
Welsh	cym_Latn	Low	Baseline	0	33.13	58.6
			RESETOX _{ETOX}	0	33.16	58.62
Danish	dan_Latn	High	Baseline	221	40.78	65.41
			RESETOX _{ETOX}	85	40.5*	65.19*
German	deu_Latn	High	Baseline	191	34.91	62.2
			RESETOX _{ETOX}	71	34.89	62.13
Southwestern Dinka	dik_Latn	Low	Baseline	25725	3.51	21.13
			RESETOX _{ETOX}	11737	3.51	21.06
Dyula	dyu_Latn	Low	Baseline	2009	1.65	19.19
			RESETOX _{ETOX}	1263	1.63	19.18

Language	Code	Resource	Model	Holistic Bias	FLORES 200	
				ETOX	BLEU	CHRF
Esperanto	epo_Latn	Low	Baseline	0	32.96	61.85
			RESETOX _{ETOX}	0	32.86	61.84
Estonian	est_Latn	High	Baseline	1027	19.49	53.27
			RESETOX _{ETOX}	622	19.45	53.23
Basque	eus_Latn	High	Baseline	4377	14.77	52.97
			RESETOX _{ETOX}	745	14.68	52.8*
Ewe	ewe_Latn	Low	Baseline	7012	11.76	38.0
			RESETOX _{ETOX}	2820	11.31*	37.47*
Faroese	fao_Latn	Low	Baseline	377	20.57	45.91
			RESETOX _{ETOX}	142	20.58	45.87
Fijian	fij_Latn	Low	Baseline	3754	17.68	46.24
			RESETOX _{ETOX}	1633	17.59	46.13
Finnish	fin_Latn	High	Baseline	1935	18.93	53.08
			RESETOX _{ETOX}	1348	18.93	53.05
Fon	fon_Latn	Low	Baseline	8580	2.49	18.68
			RESETOX _{ETOX}	4195	2.48	18.85
Friulian	fur_Latn	Low	Baseline	409	28.01	54.7
			RESETOX _{ETOX}	115	27.52*	54.31*
Nigerian Fulfulde	fuv_Latn	Low	Baseline	347	1.95	20.38
			RESETOX _{ETOX}	232	1.96	20.39
West Central Oromo	gaz_Latn	Low	Baseline	10	3.52	37.28
			RESETOX _{ETOX}	2	3.52	37.28
Scottish Gaelic	gla_Latn	Low	Baseline	1416	15.42	48.04
			RESETOX _{ETOX}	462	15.4	48.01
Irish	gle_Latn	Low	Baseline	732	23.29	50.04
			RESETOX _{ETOX}	325	23.14*	49.94*
Galician	glg_Latn	Low	Baseline	420	32.09	59.24
			RESETOX _{ETOX}	50	32.03	59.24
Guarani	grn_Latn	Low	Baseline	1135	8.98	37.66
			RESETOX _{ETOX}	489	8.98	37.66
Haitian Creole	hat_Latn	Low	Baseline	291	23.22	52.22
			RESETOX _{ETOX}	68	23.19	52.2
Hausa	hau_Latn	Low	Baseline	406	23.44	51.53
			RESETOX _{ETOX}	34	23.45	51.54
Croatian	hrv_Latn	High	Baseline	577	25.0	55.16
			RESETOX _{ETOX}	388	24.94	55.08*
Ilocano	ilo_Latn	Low	Baseline	1446	23.41	53.18
			RESETOX _{ETOX}	709	23.07*	53.0
Indonesian	ind_Latn	High	Baseline	14220	43.25	68.46
			RESETOX _{ETOX}	12338	43.01*	68.16*
Icelandic	isl_Latn	High	Baseline	13	19.8	46.74
			RESETOX _{ETOX}	7	19.81	46.73

Language	Code	Resource	Model	Holistic Bias	FLORES 200	
				ETOX	BLEU	CHRf
Javanese	jav_Latn	Low	Baseline	524	26.28	55.41
			RESETOX _{ETOX}	179	26.22*	55.35*
Kabyle	kab_Latn	Low	Baseline	4	6.41	29.28
			RESETOX _{ETOX}	0	6.33	29.26
Jingpho	kac_Latn	Low	Baseline	55	11.17	37.79
			RESETOX _{ETOX}	15	11.18	37.8
Kamba	kam_Latn	Low	Baseline	0	4.46	29.44
			RESETOX _{ETOX}	0	4.43	29.41
Kabiye	kbp_Latn	Low	Baseline	0	5.64	25.6
			RESETOX _{ETOX}	0	5.64*	25.6*
Kabuverdianu	kea_Latn	Low	Baseline	57	17.54	46.42
			RESETOX _{ETOX}	9	17.57	46.36
Kikuyu	kik_Latn	Low	Baseline	538	10.58	37.56
			RESETOX _{ETOX}	127	10.49*	37.38*
Kinyarwanda	kin_Latn	Low	Baseline	1623	15.46	47.62
			RESETOX _{ETOX}	549	15.5	47.48*
Kimbundu	kmb_Latn	Low	Baseline	901	2.96	28.54
			RESETOX _{ETOX}	46	2.96	28.48
Northern Kurdish	kmr_Latn	Low	Baseline	0	10.21	39.03
			RESETOX _{ETOX}	0	10.21*	39.03*
Central Kanuri (Latin script)	knc_Latn	Low	Baseline	0	2.21	17.95
			RESETOX _{ETOX}	0	2.2	17.94
Kikongo	kon_Latn	Low	Baseline	2751	17.54	47.11
			RESETOX _{ETOX}	1903	17.54	47.1
Ligurian	lij_Latn	Low	Baseline	3	15.5	45.46
			RESETOX _{ETOX}	0	15.52	45.46
Limburgish	lim_Latn	Low	Baseline	8	10.77	44.57
			RESETOX _{ETOX}	0	10.7	44.5*
Lingala	lin_Latn	Low	Baseline	340	17.65	49.62
			RESETOX _{ETOX}	134	17.66	49.54
Lithuanian	lit_Latn	High	Baseline	390	19.67	52.06
			RESETOX _{ETOX}	224	19.67	52.05
Lombard	lmo_Latn	Low	Baseline	24	6.24	35.16
			RESETOX _{ETOX}	2	6.24	35.1
Latgalian	ltg_Latn	Low	Baseline	26	14.79	43.46
			RESETOX _{ETOX}	3	14.81	43.5
Luxembourgish	ltz_Latn	Low	Baseline	34	22.11	54.22
			RESETOX _{ETOX}	6	22.1	54.2
Luba-Kasai	lua_Latn	Low	Baseline	1234	6.31	37.64
			RESETOX _{ETOX}	317	6.07*	37.42*
Ganda	lug_Latn	Low	Baseline	246	7.26	39.31
			RESETOX _{ETOX}	16	7.25	39.3

Language	Code	Resource	Model	Holistic Bias	FLORES 200	
				ETOX	BLEU	CHRF
Luo	luo_Latn	Low	Baseline	23855	10.47	40.06
			RESETOX _{ETOX}	16351	10.24*	39.84*
Mizo	lus_Latn	Low	Baseline	2148	9.83	37.44
			RESETOX _{ETOX}	662	9.7*	37.23*
Standard Latvian	lvs_Latn	High	Baseline	889	18.32	47.96
			RESETOX _{ETOX}	113	18.25	47.88
Minangkabau (Latin script)	min_Latn	Low	Baseline	20488	18.38	50.32
			RESETOX _{ETOX}	14152	18.27*	50.24
Maltese	mlt_Latn	High	Baseline	74	24.15	63.28
			RESETOX _{ETOX}	22	24.14	63.25
Mossi	mos_Latn	Low	Baseline	820	3.48	22.57
			RESETOX _{ETOX}	210	3.5	22.65
Maori	mri_Latn	Low	Baseline	163	19.27	45.13
			RESETOX _{ETOX}	49	19.15*	45.1
Dutch	nld_Latn	High	Baseline	74	25.23	56.24
			RESETOX _{ETOX}	29	25.31	56.23
Norwegian Nynorsk	nno_Latn	Low	Baseline	54	25.04	54.61
			RESETOX _{ETOX}	19	24.9*	54.48*
Norwegian Bokmål	nob_Latn	Low	Baseline	1489	30.72	59.2
			RESETOX _{ETOX}	1222	30.64*	59.15
Northern Sotho	nso_Latn	Low	Baseline	3	22.11	51.28
			RESETOX _{ETOX}	1	22.11	51.29
Nuer	nus_Latn	Low	Baseline	51	5.41	27.52
			RESETOX _{ETOX}	5	5.41	27.54
Nyanja	nya_Latn	Low	Baseline	939	13.7	48.73
			RESETOX _{ETOX}	585	13.68	48.73
Occitan	oci_Latn	Low	Baseline	39	33.17	60.78
			RESETOX _{ETOX}	1	32.65*	60.31*
Papiamentu	pap_Latn	Low	Baseline	4019	25.56	52.82
			RESETOX _{ETOX}	2679	25.15*	52.55*
Plateau Malagasy	plt_Latn	Low	Baseline	270	16.03	52.11
			RESETOX _{ETOX}	109	15.98	52.02
Polish	pol_Latn	High	Baseline	179	18.41	48.58
			RESETOX _{ETOX}	77	18.39	48.55
Ayacucho Quechua	quy_Latn	Low	Baseline	0	2.09	27.18
			RESETOX _{ETOX}	0	2.12	27.15
Romanian	ron_Latn	High	Baseline	221	34.04	60.69
			RESETOX _{ETOX}	68	33.81*	60.47*
Rundi	run_Latn	Low	Baseline	377	11.47	43.36
			RESETOX _{ETOX}	121	11.49	43.27*
Sango	sag_Latn	Low	Baseline	5	9.06	36.0
			RESETOX _{ETOX}	1	8.95	35.87

Language	Code	Resource	Model	Holistic Bias	FLORES 200	
				ETOX	BLEU	CHRF
Sicilian	scn_Latn	Low	Baseline	14268	5.92	37.26
			RESETOX _{ETOX}	9330	5.81	37.21
Slovak	slk_Latn	High	Baseline	23	28.56	56.4
			RESETOX _{ETOX}	14	28.47	56.35
Slovenian	slv_Latn	High	Baseline	575	25.01	53.43
			RESETOX _{ETOX}	425	24.99	53.39*
Samoan	smo_Latn	Low	Baseline	2854	25.56	49.67
			RESETOX _{ETOX}	1190	25.32*	49.37*
Shona	sna_Latn	Low	Baseline	103	12.9	48.23
			RESETOX _{ETOX}	93	12.87	48.17
Somali	som_Latn	Low	Baseline	99	11.54	45.77
			RESETOX _{ETOX}	58	11.5	45.72
Southern Sotho	sot_Latn	High	Baseline	18571	18.37	48.49
			RESETOX _{ETOX}	14650	18.35	48.49
Sardinian	srd_Latn	Low	Baseline	24	25.56	54.71
			RESETOX _{ETOX}	9	25.39*	54.58*
Swati	ssw_Latn	Low	Baseline	0	9.91	47.75
			RESETOX _{ETOX}	0	9.82	47.66
Sundanese	sun_Latn	Low	Baseline	184	18.37	50.62
			RESETOX _{ETOX}	64	18.25*	50.53*
Swedish	swe_Latn	High	Baseline	333	39.62	65.13
			RESETOX _{ETOX}	88	39.8*	65.19
Swahili	swh_Latn	High	Baseline	569	32.08	60.75
			RESETOX _{ETOX}	229	32.02	60.61*
Silesian	szl_Latn	Low	Baseline	166	16.98	47.49
			RESETOX _{ETOX}	68	16.97	47.45
Tagalog	tgl_Latn	High	Baseline	446	31.37	58.08
			RESETOX _{ETOX}	299	31.27	58.07
Tok Pisin	tpi_Latn	Low	Baseline	3590	18.33	42.94
			RESETOX _{ETOX}	1419	17.09*	41.88*
Tswana	tsn_Latn	High	Baseline	11558	21.04	49.18
			RESETOX _{ETOX}	4475	20.92	49.08*
Tsonga	tso_Latn	Low	Baseline	2885	21.57	52.12
			RESETOX _{ETOX}	2117	21.56	52.1
Turkmen	tuk_Latn	Low	Baseline	556	10.69	40.33
			RESETOX _{ETOX}	377	10.52	40.32
Tumbuka	tum_Latn	Low	Baseline	1179	9.96	37.71
			RESETOX _{ETOX}	831	9.89*	37.63
Twi	twi_Latn	Low	Baseline	29683	11.2	37.27
			RESETOX _{ETOX}	7573	10.01*	35.82*
Umbundu	umb_Latn	Low	Baseline	35	2.34	30.07
			RESETOX _{ETOX}	22	2.35	30.1

Language	Code	Resource	Model	Holistic Bias	FLORES 200	
				ETOX	BLEU	CHRf
Northern Uzbek	uzn_Latn	High	Baseline	0	15.48	52.79
			RESETOX _{ETOX}	0	15.51	52.61*
Venetian	vec_Latn	Low	Baseline	1177	14.63	48.99
			RESETOX _{ETOX}	895	14.43*	48.91
Vietnamese	vie_Latn	High	Baseline	2370	38.46	56.47
			RESETOX _{ETOX}	1085	38.48	56.48
Waray	war_Latn	Low	Baseline	3734	28.59	56.11
			RESETOX _{ETOX}	2052	28.59	56.1
Wolof	wol_Latn	Low	Baseline	1	4.99	24.67
			RESETOX _{ETOX}	0	5.0	24.65
Xhosa	xho_Latn	High	Baseline	0	13.67	53.03
			RESETOX _{ETOX}	0	13.67	53.02
Yoruba	yor_Latn	Low	Baseline	18735	4.29	24.08
			RESETOX _{ETOX}	16099	4.26	24.04
Standard Malay	zsm_Latn	High	Baseline	797	37.57	65.74
			RESETOX _{ETOX}	508	37.53	65.71
Zulu	zul_Latn	High	Baseline	34	17.24	56.66
			RESETOX _{ETOX}	6	17.23	56.65
Central Atlas Tamazight	tzm_Tfng	Low	Baseline	13	5.37	28.21
			RESETOX _{ETOX}	4	5.23*	27.83*
Dzongkha	dzo_Tibt	Low	Baseline	0	0.52	39.24
			RESETOX _{ETOX}	0	0.52*	39.24*

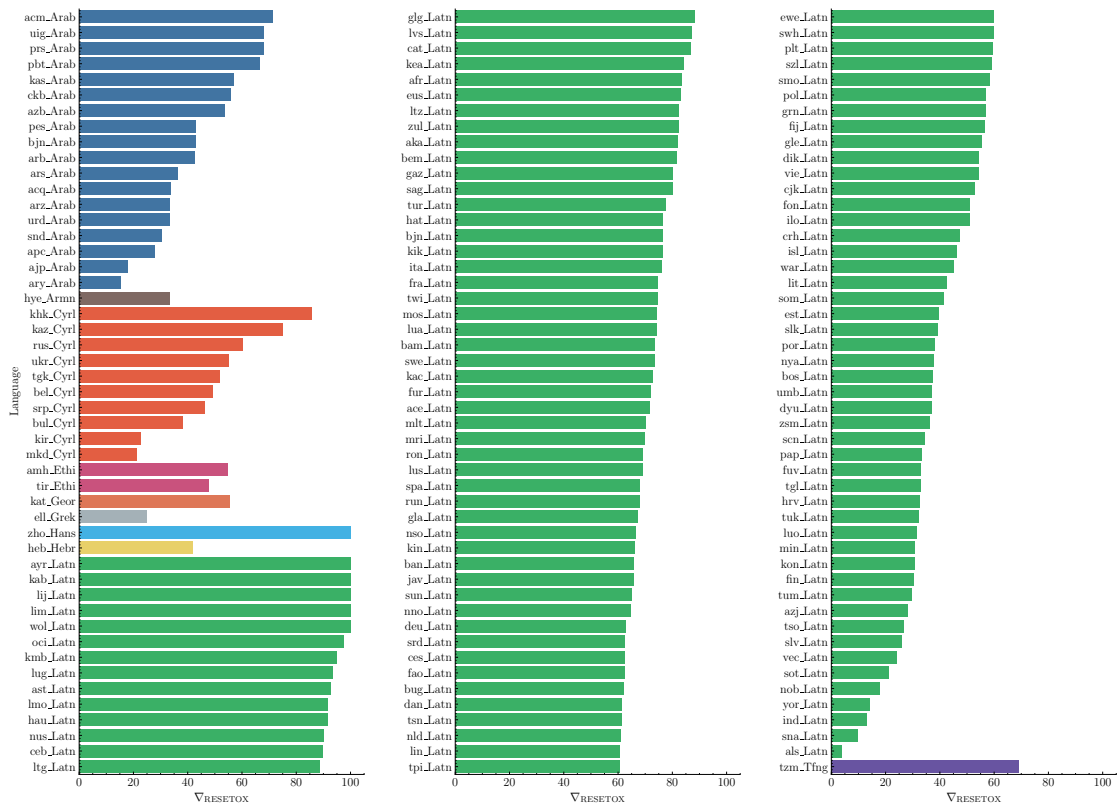


Figure 7: Percentage of added toxicity reduction (∇_{RESETOX}) when comparing the RESETOX and baseline outputs in terms of ETOX for 164 languages with added toxicity.

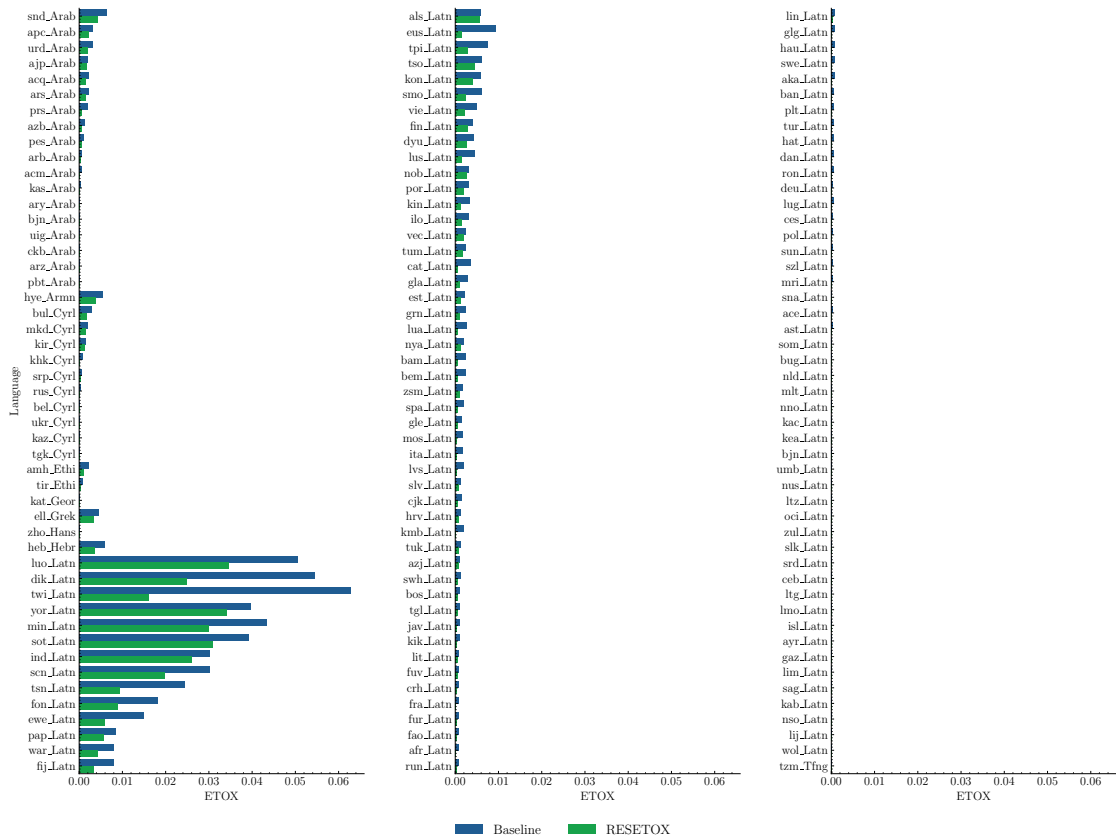


Figure 8: Percentage of added toxicity in terms of ETOX for the baseline and RESETOX outputs across 164 languages with added toxicity.

Using Machine Translation to Augment Multilingual Classification

Adam King

GumGum

aking@gumgum.com

Abstract

An all-too-present bottleneck for text classification model development is the need to annotate training data and this need is multiplied for multilingual classifiers. Fortunately, contemporary machine translation models are both easily accessible and have dependable translation quality, making it possible to translate labeled training data from one language into another. Here, we explore the effects of using machine translation to fine-tune a multilingual model for a classification task across multiple languages. We also investigate the benefits of using a novel technique, originally proposed in the field of image captioning, to account for potential negative effects of tuning models on translated data. We show that translated data are of sufficient quality to tune multilingual classifiers and that this novel loss technique is able to offer some improvement over models tuned without it.

1 Introduction

One of the most common uses of machine learning for natural language processing (NLP) is the classification of text into one of multiple mutually-inclusive or mutually-exclusive labels. Recently, generative LLMs, such as PaLM (Chung et al., 2022) and ChatGPT (Ouyang et al., 2022) have shown exciting and impressive capabilities to do zero- or few-shot prompting, classify text given only a few examples for the task across a variety of languages. Nevertheless, it is still the case that

the highest performing and most efficient means to classify text is the use of a bespoke classifier trained with hundreds or thousands labeled examples (Pires et al., 2019), particularly when the task requires a level of human-like subjectivity or general reasoning ability (Kocóń et al., 2023, see discussion). To this end, finding or creating a corpus of labeled examples is a necessary step in the creation of any classifier.

For high-resource languages like English, which have many existing labeled corpora available and large populations of annotators on crowd-sourced workers such as Amazon Mechanical Turk, the challenge of creating or finding training and evaluation data can be costly, but not prohibitively so. Yet, for lower-resourced languages which lack existing annotated corpora and have smaller or even non-existent populations on these large annotation platforms, acquiring the required training data can prove to be much more difficult. Moreover, if the model is intended to be able to perform the same classification across multiple languages, the time and effort required to annotate training data becomes multiplicative. Fortunately, classification is not alone in the applications of machine learning in NLP. Machine translation (MT) has seen major improvements in recent years (Stahlberg, 2020), accelerated by the adoption of the transformer architecture (Vaswani et al., 2017).

To date, several options for high quality machine translation currently exist, between API services and open-source models. MT API services, such as Google translate, have become nearly ubiquitous, provide high quality translations, while still being relatively inexpensive. In fact, in one experiment, translating data using Google translate into English and using existing English-trained classifier models outperformed certain models trained

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

on the original language directly (Araujo et al., 2016). In addition to MT API services, several open-source translation models are easily available, such as the multilingual M2M100 model (Fan et al., 2020), NLLB200 model (Team et al., 2022) or the over 1400 models trained by the University of Helsinki (Tiedemann and Thottingal, 2020), with many of these models have performance that approaches or exceeds that of MT APIs (Stahlberg, 2020).

With this in mind, it may be the case that translating an existing, labeled dataset with one of the aforementioned MT options is a feasible alternative to creating a novel dataset directly in that language. This has several benefits. Firstly, it avoids the problem of existing corpora or annotation options not existing for the language in question. Secondly, it minimizes the data needed for multilingual models and allows annotations for one language to serve another. Here, we ask if it is possible to use MT to train a multilingual model, given only original, annotated data for a single language.

Of course, the potential benefits of using MT to train a multilingual model are still affected by the old machine learning adage: garbage in, garbage out. Even the best translations, either human or machine, will lose some of the information of the original language, which will inevitably lead to dropped performance for a model trained on the translated examples. Fortunately, the problem of training models using semantically similar but imperfect pairs of data is not unique to the task at hand and there is a growing body of research which may provide some benefit. In particular, image captioning is a task to generate the ideal natural language text caption for an image and these captioning models must learn to represent semantically related data from very different modalities similarly, i.e., text and images (Li et al., 2021). In this way, image captioning is somewhat analogous to the task of training on translated data, where we want to have semantically identical text from different languages predicted to have the same labels. As a result, we ask in addition whether some of the model training techniques used in image captioning models can lead to improved performance for multilingual models trained using MT data.

2 Related Work

This work is by no means the first to suggest the usage of machine translation to create or augment

datasets for lower resourced languages. Wei and Pal (2010) and Pan et al. (2011) augmented Chinese language corpora with annotated data translated from English to improve the performance of a Chinese-language sentiment analysis model. On the other hand, Barriere and Balahur (2020) and Ghafoor et al. (2021) used existing API translation services to translate annotated data from English into lower-resourced languages and trained classifiers solely on these translated data, finding that classifiers trained on translated data were fairly accurate but did see drops in performance, likely due to the effects of imperfect translations of the training data.

It should be noted that training a model from scratch is not the only means to create an accurate classifier, particularly for lower-resourced languages. Large multilingual transformer models such as m-BERT (Devlin et al., 2018), XLM-ROBERTA (Conneau et al., 2019) or GPT-3 (Brown et al., 2020) have been shown to have the ability to generalize from one language to the other, i.e., train in one language and improve test performance in another language, (Pires et al., 2019), but benefits of this vary on the languages in question, with languages that share closer genealogical origin or structural similarities benefiting more from inter-language transfer. Regardless, training a model with examples of a particular language dependably yields the best classifier for new data in that language.

Nevertheless, to date there has been no investigation of how fine-tuning large multilingual transformer models on translated data affects final performance compared to simple interlanguage transfer. Moreover, previous work to train models using translated data employed a naive approach, treating translated data as if it were no different than original, untranslated data which annotated itself. In this work, we investigate both how multilingual transformer models trained on translated data perform compared to interlanguage transfer and explore a means to mitigate imperfect translation quality when creating these training datasets.

3 Image captioning and Image-Text Contrastive Loss

Image-text Contrastive (ITC) loss is a technique used when training multimodal models to caption images with natural language descriptions (Li et al., 2021). For example, BLIP (Li et al., 2022)

is a image-captioning model that was trained with a mix of human- and artificially-annotated images where ITC loss was integral to the models ability to learn from noisy, artificially-annotated data. ITC loss, then, has been shown to mitigate negative effects of both noise and different modalities for multimodal models.

At an intuitional level, these captioning models decompose text and images into a shared embedding space and ITC loss seeks to penalize cases where related image-text pairs are dissimilar in this shared embedding space. In other words, ITC looks seeks to bring semantically related items from disparate modalities closer in a shared embedded space and has empirically improved image-captioning models, with little impact on training time or resources.

Training multilingual classification models with translated data bears a similarity to captioning, though rather than have semantically related examples from different modalities, there are semantically parallel data in different languages. That being the case, we will be a slightly modified form of ITC loss, namely original-translated contrastive (OTC) loss, to enforce similarity within a batch between data from the original language and its translated counterpart. Like ITC loss, OTC loss penalizes a transformer model for dissimilar embedding representations for translated pairs. One way to think of it is that this loss encourages the model to embed sentences with the same meaning identically, regardless of language.

In detail, we implement OTC loss as follows. We begin by deriving a probability of each original/translated pairing in a training minibatch, p^{o2t} and p^{t2o} , that is, which original examples pairs with which translated example and vice versa.

$$p_m^{o2t} = \frac{\exp(s(O, T_m)/\tau)}{\sum_m^M \exp(s(O, T_m)/\tau)} \quad (1)$$

$$p_m^{t2o} = \frac{\exp(s(T, O_m)/\tau)}{\sum_m^M \exp(s(T, O_m)/\tau)} \quad (2)$$

Here, $s(T, O)$ is a similarity function between the original, untranslated data and the translated examples in a minibatch. We compute $s(T, O)$ by first extracting and normalizing the embedding for the initial [CLS] token after the final attention head of the encoder stack in M-BERT, computing a pairwise dot product for all possible pairs of original and translated data and dividing by τ , which is a learnable parameter. We then apply the

softmax function as a way to represent the likelihood of each original/translated match. Ideally, each correct original/translated pair will have the most similar embeddings, resulting in a value close to 1 after softmax. As a final step, we compute the cross-entropy between the result of the previous step and a target vector which encodes the correct original/translated pairs, weighting this by a hyperparameter, α_{otc} . Following BLIP (Li et al., 2022), we set $\alpha_{otc} = .4$ for all runs.

$$\ell_{otc} = \alpha_{otc} * \frac{1}{2} \mathbb{E}_{(O,T)} [H(\mathbf{y}^{o2t}(O), \mathbf{p}^{o2t}(O)) + H(\mathbf{y}^{t2o}(T), \mathbf{p}^{t2o}(T))] \quad (3)$$

4 Experiments

4.1 Data

For these experiments, we use a multilingual dataset of Amazon product reviews across 6 languages: English, Spanish, French, German, Chinese and Japanese (Keung et al., 2020). This dataset is comprised of over 1 million total examples, split into a train and test partition. The reviews are equally distributed across the six languages, as well as the total stars given to the reviewed product (1-5) for both the train and test partition, i.e., each number of stars comprises 20% of the examples for that language. This dataset is particularly useful due to its size, number of available languages and presence of an established training and test data split.

We began by translating each review from the training partition of the original dataset into each of the other respective languages and assigned the same star value to the review (see example 1), i.e., if a review was originally in English and had star star, when translating it into French it would also be labeled with one star. We did this translation once before carrying out the rest of the experiment to ensure each classifier would be trained on the same set of translations. To translate, we used a single multilingual translation model, M2M100 (Fan et al., 2020). We chose to use a single multilingual translation model in order to mitigate any potential differences from translation quality coming from different machine translation architectures.

4.2 Experiment design

To investigate any potential improvement in classifier accuracy with the use OTC loss, we fine-tuned

id	translated	language	text	stars
1	0	en	My daughter really likes the backpack and ...	5
1	1	es	Mi hija realmente le gusta el bolsillo y ...	5
...
2	0	en	This product is BS, I washed my face with hot water ...	1
2	1	fr	Ce produit est BS, je me suis lavé le visage à l'eau chaude ...	1
...

Figure 1: Example original and translated data. Each unique review (id) in the original dataset was translated to the other languages and assigned the same star value. Texts truncated here for formatting.

pretrained transformer models on datasets that included original, untranslated data for a single language¹ and only translated data for all others in the six language set. As an example, in one training run, the model would be tuned on the original English training data and only translated data for all other languages, which were translated from the set of the original English data. We did this for all six languages in the original set to ensure any results were not restricted to one language in the dataset. Though the exact training examples varied for each model, we tested each on the original testing split of the dataset, which was solely comprised of original data, i.e., non-translated, for the six languages.

In each case, we tuned a multilingual DISTIL-BERT model (Sanh et al., 2019), a distilled version of the original multilingual M-BERT (Devlin et al., 2018), to predict the number of stars on a review as a categorical classification problem, using categorical cross-entropy loss and varying between using OTC loss as an additional loss parameter between runs. We chose to use a distilled variant of BERT due to the distilled variants increased speed of training, while still maintaining 97% of overall language understanding of the original.

Because of the mechanics of OTC loss, each translated datum must have an original match in the minibatch and each original must have at least one translated variant. As such, we constructed minibatches during training such that half the samples were always original, untranslated data and the other half were a randomly selected translated example for each original datum. For each original example, we randomly selected a translated example from the other languages, meaning that the model saw an equal number of original and

translated examples during tuning overall, though it saw far fewer individual examples of each translated language, i.e., roughly $\frac{1}{5}$. For simplicity, we restricted our tests to a 1:1 original:translated ratio and we used the same batch sampling method for runs without OTC loss, to make results more easily comparable.

For each tuning run, we used a batch size of 32 (16 original and 16 translated examples per batch)² and used the AdamW (Loshchilov and Hutter, 2017) optimizer with a linear warm-up of 500 updates with a learning rate of $2e-5$. All training was done on G5.2XLARGE AWS instances which contain NVIDIA A10G GPUs. We tuned 3 separate tuning runs for each set of hyperparameters and report their mean values in the next section.

5 Results

In these experiments, we asked two simple questions: 1) how feasible is it to tune a multilingual transformer model on translated data and 2) does the inclusion of OTC loss improve model performance for languages where only translated training data was used.

In answer to the first, for each of the six languages in the original dataset, models fine-tuned with translated data showed higher F1-micro scores³ on the held-out test set, compared to models trained with only original data for a single language (see Table 1). As was expected from Pires et al. (2019), even if a model was never exposed to data for a language, original or translated, the final model did have F1-micro greater than chance for that language (which would be 20% for a balanced, 5-label problem), indicating

¹We restricted the experimental conditions to only including a single language’s original data, rather than use the full set of $6! = 720$ possible permutations of language combinations for the sake of efficiency and resources.

²For baseline conditions where there was no translated data, mini-batching happened as normal with 32 examples original, untranslated data per batch.

³F1-micro is an example-weighted version of the F1-score, which is the harmonic mean of precision and recall. For more details on F1-score, see (Jurafsky and Martin, 2008).

Language	F1-micro		
	No data	Translated	Original
EN	0.407	0.481	0.554
FR	0.379	0.468	0.544
DE	0.359	0.465	0.581
ES	0.376	0.474	0.55
JA	0.307	0.396	0.543
ZH	0.352	0.372	0.458

Table 1: F1-micro for models trained with no samples for the specified language (No data), with only translated samples (Translated) and with the original training data for that language (Original). All languages saw a sizeable boost to performance over their respective baselines when using translated data (.02-.11) but all languages did perform markedly better when given actual data for each language.

there was interlingual knowledge transfer happening within the model during training. Moreover, it appears that there was more transfer between related, similar languages, compared to more dissimilar languages; models trained with data for a European language showed higher performance on other European languages, compared to Japanese or Chinese. Nevertheless, for all languages, the use of translated data did show a noticeable improvement (.02-.11), though for each language, models trained with only translated data did underperform models trained with the full set of original, untranslated training examples for that language (.07-.12).

That said, it is clear that the use of translated training data does improve model performance, even if the trained model only sees translated examples for that language. It should also be noted that due to the batching and sampling strategy used here, models trained with translated data saw far fewer examples of each language where they only saw translated data. That is, because each original review was paired with a single translated example out of five possible translated, these models were exposed to roughly one fifth of the data for translated languages and still saw a sizable boost in performance.

Moving on to the effect of OTC loss, Table 3 shows the mean F1-micro per language in the testing set, for models fine-tuned using original data for the specified language and translations for all other languages. For all languages, models trained using OTC loss saw an improvement over models trained without for all languages except Chinese, which showed a mixed set of negligible differences or lowered performance. However,

Language	F1-micro	
	No OTC	OTC
EN	0.479	0.483
FR	0.464	0.472
DE	0.463	0.467
ES	0.472	0.476
JA	0.393	0.399
ZH	0.368	0.376

Table 2: Comparison on final performance per language for models that only included translated examples for the specified language. Though the gain was less than .1, each language consistently performed better when trained with OTC loss.

these values include runs where the specific language was included as original, untranslated data. When averaging across all runs where a language in the testing set was only represented by translated data, OTC loss shows an improvement over models trained without it for all languages. Table 2 shows the mean F1-micro for all models trained where the specified language was not the original language.

To ensure that the results here were in fact statistically significant, we fit a linear mixed-effect model to predict final model F1-micro for a language, given the hyperparameters of a particular tuning run. Mixed-effect models are able to accurately evaluate the contribution of different fixed-effect independent variables, e.g., whether OTC was used when training a particular model, on dependent variables, e.g., the final accuracy of the trained model, all the while being robust to expected random variance between trials, e.g., because of random initialization and batching, some deep learning models score higher than others with identical hyperparameters (see Baayen et al. (2008), Jaeger (2008) for more).

This statistical model was fit to predict per-language test f1-micro, given a random effect of each model run and three fixed effects: i) the tested language, ii) the identity of the single original language and iii) whether OTC loss was added. OTC was found to have a significant, positive effect (COEF=0.036, STD.ERROR=0.017, for all model details see 2), indicating that even after taking into consideration differences between languages and random variance for each multilingual model, the inclusion of OTC loss did yield an improved final model F1-micro.

Orig. Training Language	OTC	EN	FR	DE	ES	JA	ZH
EN	No OTC	0.548	0.488	0.493	0.489	0.425	0.423
	OTC	0.553	0.507	0.522	0.512	0.434	0.422
FR	No OTC	0.504	0.539	0.504	0.493	0.424	0.426
	OTC	0.512	0.539	0.517	0.511	0.428	0.412
DE	No OTC	0.514	0.495	0.577	0.495	0.436	0.427
	OTC	0.524	0.506	0.581	0.506	0.449	0.425
ES	No OTC	0.506	0.497	0.500	0.544	0.433	0.419
	OTC	0.523	0.510	0.518	0.548	0.441	0.413
JA	No OTC	0.470	0.460	0.477	0.468	0.526	0.436
	OTC	0.493	0.474	0.499	0.487	0.522	0.424
ZH	No OTC	0.486	0.439	0.441	0.444	0.398	0.482
	OTC	0.488	0.467	0.473	0.472	0.421	0.503

Table 3: F1-micro results on untranslated test data. Each row shows the per-language performance for models trained with original data for the specified language and translated data for all other languages, using OTC loss and without. Each cell shows the mean of 3 runs per condition. Bolded values show a difference of .03 or greater.

6 Discussion and future directions

We investigated the feasibility of using translated text to fine-tune a multilingual transformer model, as well as any potential gains by utilizing a novel application of deep learning technique to improve performance. We found that models trained using only translated data for a language do show a noticeable improvement over baselines, though as expected, there was still a performance drop from using original, untranslated data for that language. We also found that slight further gains can be achieved by the use of OTC loss, suggesting that training the model in such a way where it is sensitive to potential data issues improves its ability to generalize.

Granted, this is a very open problem and results of using translated data to tune a multilingual classifier will vary highly depending on the quality of MT model used, architecture of the classifier being tuned and the type of classification being modeled. Nevertheless, the results here are exciting for multiple reasons. Firstly, as suggested by previous works (Shalunts et al., 2016, as an example), MT is useful tool for language-specific dataset creation when creating a dataset for that language directly may prove difficult. In this case, we showed that M-BERT models tuned on translated examples showed large gains over simple multilingual transfer during training. This is particularly interesting given that for each translated language, the model was only given a fraction of samples compared to the original language due to the 1:1 ratio of original and translated data. A future direction

for this work may be to adjust this ratio or the number of languages in the dataset to investigate how this affects model training. Secondly, the use of OTC loss was shown to lead to a small, but robust boost to performance. This suggests that methods of mitigating the natural effects of translation have a potential to bridge the gap, so to speak, between models trained on translated data and on datasets in the target language directly. Particularly relevant, Chinese, which is linguistically dissimilar from the majority of languages in the set used here, showed a mixed ability to benefit from training with other languages, but a clearer improvement using OTC loss. This may suggest that OTC loss is able to mitigate structural differences between languages and a future direction for this may be to explore exactly how OTC loss affects individual examples and how other noise-reduction techniques may lead to further gains in model performance.

Putting this together, this is an indication that MT-augmented datasets stand as a good first step for developing multilingual classification models. Given that MT can quickly and efficiently expand an annotated dataset from one language into another and that translated dataset is of sufficient quality to improve over basic interlingual transfer, this technique has great potential to expanding classification tasks to new languages quickly. In addition, OTC loss may be able to slightly but significantly increase the quality of these models with no additional data. All in all, we are confident that the use of MT augmentation is an exciting and interesting topic for future exploration.

7 Acknowledgements

This work was carried out as a part of the R&D for GumGum’s Verity product. Special thanks go to the members of the AI team (names, names, names) for their suggestions and audience during brainstorming, development and analysis of this project.

References

- Araujo, Matheus, Julio Reis, Adriano Pereira, and Fabricio Benevenuto. 2016. An evaluation of machine translation for multilingual sentence-level sentiment analysis. In *Proceedings of the 31st annual ACM symposium on applied computing*, pages 1140–1145.
- Baayen, R Harald, Douglas J Davidson, and Douglas M Bates. 2008. Mixed-effects modeling with crossed random effects for subjects and items. *Journal of memory and language*, 59(4):390–412.
- Barriere, Valentin and Alexandra Balahur. 2020. Improving sentiment analysis over non-english tweets using multilingual transformers and automatic translation for data-augmentation. *arXiv preprint arXiv:2010.03486*.
- Brown, Tom B., Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. *CoRR*, abs/2005.14165.
- Chung, Hyung Won, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Conneau, Alexis, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.
- Fan, Angela, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation. *CoRR*, abs/2010.11125.
- Ghafoor, Abdul, Ali Shariq Imran, Sher Muhammad Daudpota, Zenun Kastrati, Rakhi Batra, Mudasir Ahmad Wani, et al. 2021. The impact of translating resource-rich datasets to low-resource languages through multi-lingual text processing. *IEEE Access*, 9:124478–124490.

- Jaeger, T Florian. 2008. Categorical data analysis: Away from anovas (transformation or not) and towards logit mixed models. *Journal of memory and language*, 59(4):434–446.
- Jurafsky, Daniel and James H Martin. 2008. Speech and language processing: An introduction to speech recognition, computational linguistics and natural language processing. *Upper Saddle River, NJ: Prentice Hall*.
- Keung, Phillip, Yichao Lu, György Szarvas, and Noah A. Smith. 2020. The multilingual amazon reviews corpus. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*.
- Kocoń, Jan, Igor Cichecki, Oliwier Kaszyca, Mateusz Kochanek, Dominika Szydło, Joanna Baran, Julita Bielaniec, Marcin Gruza, Arkadiusz Janz, Kamil Kanclerz, et al. 2023. Chatgpt: Jack of all trades, master of none. *Information Fusion*, page 101861.
- Li, Junnan, Ramprasaath R. Selvaraju, Akhilesh Deepak Gotmare, Shafiq R. Joty, Caiming Xiong, and Steven C. H. Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. *CoRR*, abs/2107.07651.
- Li, Junnan, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation.
- Loshchilov, Ilya and Frank Hutter. 2017. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101.
- Ouyang, Long, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Pan, Junfeng, Gui-Rong Xue, Yong Yu, and Yang Wang. 2011. Cross-lingual sentiment classification via bi-view non-negative matrix tri-factorization. In *Advances in Knowledge Discovery and Data Mining: 15th Pacific-Asia Conference, PAKDD 2011, Shenzhen, China, May 24-27, 2011, Proceedings, Part I 15*, pages 289–300. Springer.
- Pires, Telmo, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual bert? *CoRR*, abs/1906.01502.
- Sanh, Victor, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108.
- Shalunts, Gayane, Gerhard Backfried, and Nicolas Commeignes. 2016. The impact of machine translation on sentiment analysis. *Data Analytics*, 63:51–56.
- Stahlberg, Felix. 2020. Neural machine translation: A review. *Journal of Artificial Intelligence Research*, 69:343–418.
- Team, NLLB, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.
- Tiedemann, Jörg and Santhosh Thottingal. 2020. OPUS-MT – building open translation services for the world. In Martins, André, Helena Moniz, Sara Fumega, Bruno Martins, Fernando Batista, Luisa Coheur, Carla Parra, Isabel Trancoso, Marco Turchi, Arianna Bisazza, Joss Moorkens, Ana Guerberof, Mary Nurminen, Lena Marg, and Mikel L. Forcada, editors, *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 479–480, Lisboa, Portugal, November. European Association for Machine Translation.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Wei, Bin and Christopher Pal. 2010. Cross lingual adaptation: an experiment on sentiment classifications. In *Proceedings of the ACL 2010 conference short papers*, pages 258–262.

Mixed Linear Model Regression Results

```

=====
Model:                MixedLM   Dependent Variable:  test_acc
No. Observations:    108       Method:              REML
No. Groups:          18        Scale:               0.0038
Min. group size:     6         Log-Likelihood:     111.7634
Max. group size:     6         Converged:           Yes
Mean group size:     6.0
-----

```

	Coef.	Std.Err.	z	P> z	[0.025	0.975]
Intercept	0.465	0.024	19.128	0.000	0.418	0.513
otc[T.True]	0.036	0.017	2.105	0.035	0.002	0.069
original_lang[T.en]	-0.020	0.028	-0.714	0.475	-0.074	0.035
original_lang[T.es]	-0.012	0.028	-0.432	0.666	-0.066	0.042
original_lang[T.fr]	-0.015	0.028	-0.555	0.579	-0.070	0.039
original_lang[T.ja]	-0.030	0.028	-1.093	0.274	-0.085	0.024
original_lang[T.zh]	-0.050	0.028	-1.801	0.072	-0.104	0.004
test_lang[T.en]	0.016	0.020	0.762	0.446	-0.025	0.056
test_lang[T.es]	0.003	0.020	0.130	0.897	-0.037	0.043
test_lang[T.fr]	-0.000	0.020	-0.005	0.996	-0.040	0.040
test_lang[T.ja]	-0.061	0.020	-2.972	0.003	-0.101	-0.021
test_lang[T.zh]	-0.068	0.020	-3.336	0.001	-0.108	-0.028
Group Var	0.001	0.009				

```

=====

```

Figure 2: Full model details for MLE model trained to predict F1-micro per language. OTC has a positive contribution to an increase F1-micro score, even when controlling for variance between languages and model runs.

Recovery Should Never Deviate from Ground Truth: Mitigating Exposure Bias in Neural Machine Translation

Jianfei He¹, Shichao Sun², Xiaohua Jia¹, Wenjie Li²

¹ City University of Hong Kong, Tat Chee Avenue, Kowloon, Hong Kong

² The Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong

jianfeihe-2c@my.cityu.edu.hk, bruce.sun@connect.polyu.hk

csjia@cityu.edu.hk, wenjie.li@polyu.edu.hk

Abstract

In Neural Machine Translation, models are often trained with teacher forcing and suffer from exposure bias due to the discrepancy between training and inference. Current token-level solutions, such as scheduled sampling, aim to maximize the model’s capability to recover from errors. Their loss functions have a side effect: a sequence with errors may have a larger probability than the ground truth. The consequence is that the generated sequences may deviate from the ground truth. This side effect is verified in our experiments. To address this issue, we propose using token-level contrastive learning to coordinate three training objectives: the usual MLE objective, an objective for recovery from errors, and a new objective to explicitly constrain the recovery in a scope that does not impact the ground truth. Our empirical analysis shows that this method effectively achieves these objectives in training and reduces the frequency with which the third objective is violated. Experiments on three language pairs (German-English, Russian-English, and English-Russian) show that our method outperforms the vanilla Transformer and other methods addressing the exposure bias.

1 Introduction

Like many other text generation tasks, models for Neural Machine Translation (NMT) (Bahdanau et

al., 2014) are usually trained with *teacher forcing*. During training, ground truth tokens are used as target prefixes to the decoder, and the model learns to predict the next token conditioned on the ground truth. There is a discrepancy between this training method and inference. In inference, the ground truth tokens are not available. The target prefixes to the decoder are tokens previously generated by the model, which may include some errors. This discrepancy is referred to as *exposure bias* (Bengio et al., 2015; Ranzato et al., 2016). The main concern about exposure bias is *error accumulation*. If one error happens at one step, it is incorporated into the future steps and leads to more errors. Although there are still some doubts about whether exposure bias is a big issue for text generation (He et al., 2021), more research shows that this issue matters for NMT (Wu et al., 2018; Wang and Senrich, 2020; Korakakis and Vlachos, 2022).

There are two approaches to mitigate the exposure bias, working at the token and sequence levels, respectively.

The token-level solutions, for example, *scheduled sampling* (Bengio et al., 2015; Mihaylova and Martins, 2019; Liu et al., 2021), usually use the tokens sampled from the model to replace the ground truth in training. The objective is to simulate the possible errors in inference and recover from these errors to reduce the error accumulation.

The sequence-level solutions directly maximize the total quality of the generated sequences with a sequence-level loss function (Ranzato et al., 2016; Shen et al., 2016; Edunov et al., 2018). There is still debate whether these solutions are stable and effective (Choshen et al., 2019; Kiegeland and Kreutzer, 2021).

This paper focuses on mitigating the exposure bias with token-level objectives.

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

The loss functions used in most token-level solutions have a side effect. They aim to increase the model’s capability to recover from errors by maximizing the probability of the next token conditioned on some error tokens. Consequently, a sequence with errors may have a larger probability than the ground truth, and the generated sequences may deviate from the ground truth. This side effect is verified in our experiments. We discover a missing objective behind this side effect that can explicitly constrain the recovery in a scope that does not impact the ground truth. We propose to use *token-level contrastive learning* and coordinate three training objectives: the usual Maximum Likelihood Estimation (MLE) objective, an objective for recovery from errors, and a new objective constraining the recovery. Our empirical analysis shows that this method effectively meets three objectives in training. Particularly our method reduces the frequency that the third objective is violated. We conduct experiments on German-English (De–En), Russian-English (Ru–En), and English-Russian (En–Ru). Results show that our method outperforms the vanilla Transformer and other methods addressing the exposure bias.

2 Related Work

2.1 Exposure Bias and Methods to Mitigate It

The existence of exposure bias is well recognized (Bengio et al., 2015; Ranzato et al., 2016), but its impact is still under debate. He et al. (2021) find that the distortion from exposure bias is limited in open-ended generation tasks. They hypothesize that the self-recovery ability of the language model is countering that distortion. In NMT, Wu et al. (2018) and Korakakis and Vlachos (2022) prove the *error accumulation* from exposure bias using *prefix switching*. They use different types of prefixes on the target side and measure the difference in the quality of the predictions. Typical prefixes include ground truth, predictions from the system, and random tokens. Wang and Sennrich (2020) provide indirect evidence for exposure bias in NMT. They train models with Minimum Risk Training (MRT), which has a sequence-level objective and inherently avoids exposure bias. The better performance of MRT than MLE justifies that exposure bias is harmful. Besides NMT, Chiang and Chen (2021) and Arora et al. (2022) quantify exposure bias in open-ended text generation tasks such as text completion.

Two categories of approaches have been proposed to mitigate exposure bias.

The token-level approach usually uses the tokens sampled from the model to replace the ground truth in training. Bengio et al. (2015) propose *Scheduled Sampling (SS)*, which dynamically takes samples from the model’s predictions and replaces the ground truth used for the decoder. Zhang et al. (2019) further extend the sample space with beam search and choose the candidate translation with a sentence-level metric such as BLEU. Mihaylova and Martins (2019) implement SS to Transformer (Vaswani et al., 2017) using *two-pass decoding*. The first pass gets the predictions from the model, which are used as input to the second decoder according to the scheduler. Liu et al. (2021) propose *Confidence-Aware Scheduled Sampling (CASS)* which also uses the two-pass decoding. They improve the performance by choosing the inputs to the second decoder based on the log probability of the ground truth token. Model predictions are only used when the model is confident and has a high probability (above 0.9 in their paper). Goodman et al. (2020) propose *TeaForN* to mitigate exposure bias. They use a stack of decoders to allow the model to update based on N prediction steps. Each decoder’s output is used to calculate the loss component at this decoder and is also used as the input of the next decoder. The overall loss is the weighted sum of losses from all decoders.

There are some doubts about SS. Huszár (2015) proves that SS has an improper training objective. Experiments in Mihaylova and Martins (2019) show that SS performs worse than teacher forcing for De–En. Korakakis and Vlachos (2022) use the ground truth tokens as prefixes for the decoding on a model trained with SS and find that its performance is bad compared to the MLE model. They conclude that *finetuning* with SS results in *catastrophic forgetting* (French, 1999). To avoid forgetting, they use Elastic Weight Consolidation (EWC) to regularize conditioning with model-generated prefixes. This method is similar to TFN for using a weight for prediction. But EWC works at the training parameters level, not at the loss level like TFN.

The sequence-level approach uses a sequence-level loss function and directly maximizes the total quality of the generated sequences. Ranzato et al. (2016) propose MIXER, based on

a reinforcement-learning algorithm REINFORCE. MRT (Shen et al., 2016; Wang and Sennrich, 2020) aims to minimize the risk by preference to the candidate with the largest similarity to other candidates. Edunov et al. (2018) provide a summary of classic sequence-level loss functions. There is some debate on the effectiveness of these methods. Choshen et al. (2019) identify multiple weaknesses of MIXER and MRT and suspect that they do not optimize the expected reward, while Kiegele and Kreutzer (2021) provide empirical counter-evidence to these claims.

The sequence-level approach is usually hard to converge from randomly initialized parameters and requires a baseline model trained at the token level as a starting point. In this sense, a token-level solution can be complementary to the sequence-level approach.

2.2 Using Contrastive Learning (CL) in NLP

Sun and Li (2021) apply CL to mitigate exposure bias for text summarization. They use the gold references and low-quality predictions as the positive and negative samples, respectively. The average log probability of sequences is used for the loss. Liu et al. (2022) use CL to calibrate the model. The objective is that higher-quality candidates tend to have higher log probability and are more likely to be chosen from the n-best list at the decision phase. All these methods use CL in sequence-level objectives, while our method works at the token level.

Yang et al. (2019) and Pan et al. (2021) apply CL to NMT, but they address specific issues, namely word omission errors and interim presentation for many-to-many multilingual NMT, respectively. Su et al. (2022) use CL to calibrate the model’s representation space for tokens, mitigating the issue of anisotropic distribution of token representations.

3 Approach

3.1 Discover the Missing Objective

We analyze the objectives used by the current token-level methods and discover a missing objective.

We use X and y_i to denote the source sentence and the ground truth token for step i . \hat{y}_i is a target token different from y_i at step i .

At step i , the MLE training with teacher forcing maximizes $p(y_i|X, y_1, \dots, y_{i-1})$. If the model is effectively trained, it implies that, for any \hat{y}_i ,

$$p(y_i|X, y_{<i}) > p(\hat{y}_i|X, y_{<i}). \quad (1)$$

The popular token-level methods addressing exposure bias, such as Scheduled Sampling, usually aim to enhance recovery capability from errors by maximizing $p(y_i|X, y_{<i-1}, \hat{y}_{i-1})$, which implies that, for the sampled \hat{y}_{i-1} and any \hat{y}_i ,

$$p(y_i|X, y_{<i-1}, \hat{y}_{i-1}) > p(\hat{y}_i|X, y_{<i-1}, \hat{y}_{i-1}). \quad (2)$$

Note: when \hat{y}_{i-1} is the first error, $y_{<i-1}$ are all ground truth tokens. Otherwise, $y_{<i-1}$ may include sample tokens.

However, maximizing $p(y_i|X, y_{<i-1}, \hat{y}_{i-1})$ has a side effect. Although it is good for recovery, it may impact the ground truth. If $p(y_i|X, y_{<i-1}, \hat{y}_{i-1})$ exceeds $p(y_i|X, y_{<i})$, the sequence $(y_{<i-1}, \hat{y}_{i-1}, y_i)$ may have a larger probability than the ground truth $(y_{<i-1}, y_{i-1}, y_i)$. This side effect is observed in our experiments (Subsection 5.2).

This side effect implies that the model’s prediction may deviate from the ground truth and generate a sequence with an error. This is particularly probable when beam search is used for decoding, where several \hat{y}_{i-1} tokens have a chance to remain in the hypothesis set and enter the next step during decoding.

The objective in Inequality (3) is *missing* in current training objectives:

$$p(y_i|X, y_{<i}) > p(y_i|X, y_{<i-1}, \hat{y}_{i-1}). \quad (3)$$

With this objective, the recovery is explicitly constrained in a scope not to impact the ground truth. We propose to include it in training.

These three inequalities represent three objectives that we want to achieve. We denote them as Obj_{MLE} , Obj_{Rec} and Obj_{CRec} for Inequality (1), (2), and (3), respectively. $CRec$ stands for *Constraining the Recovery*.

3.2 Token-Level Contrastive Learning

The key component in the loss function of contrastive learning is a *max* function:

$$\max\{0, \rho + S_{negative} - S_{positive}\}, \quad (4)$$

where $S_{negative}$ and $S_{positive}$ are scores for negative and positive samples, ρ is a hyperparameter for the margin. This function implies that when the score of the negative sample plus a margin is larger than the score of the positive sample, it outputs a positive loss. Otherwise, the loss is zero. The objective is that the score of the negative sample is

constrained to be at least one margin lower than the score of the positive sample.

We apply contrastive learning at the token level. The left terms in Inequality (2) and (3) are used as the scores of positive samples, while their right terms are the scores of negative samples.

3.3 Coordinate Three Objectives in One Loss Function

Three objectives in Subsection 3.1 are combined in our loss function using *multi-task learning*.

We follow the two-pass decoding in Mihaylova and Martins (2019) and Liu et al. (2021), as illustrated in Figure 1.

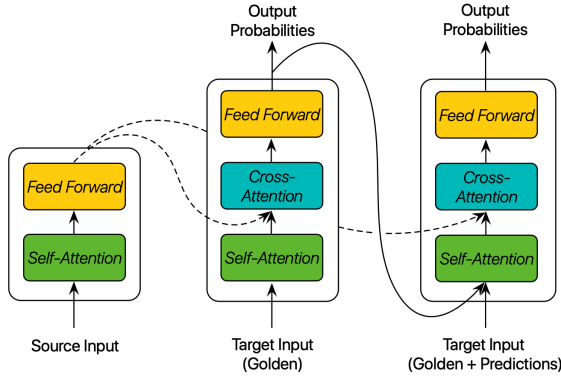


Figure 1: Scheduled sampling for the transformer with two-pass decoding (Mihaylova and Martins, 2019; Liu et al., 2021)

The first decoder is trained with teacher forcing, and its output is used for the Obj_{MLE} (Inequality 1). The Negative Log-Likelihood (NLL) with Label Smoothing (Edunov et al., 2018) is used:

$$\mathcal{L}_{MLE} = - \sum_{i=1}^n \log p(y_i | X, y_{<i}) - D_{KL}(f \parallel p(y_i | X, y_{<i})), \quad (5)$$

where f is uniform prior distribution over all tokens in the vocabulary with the size of V , $f = \frac{1}{V}$.

We use the same strategy and hyperparameters in *Confidence-Aware Scheduled Sampling* (Liu et al., 2021) to decide the inputs to the second decoder. Predicted tokens and random tokens are used as target inputs for high-confidence positions, and the ground-truth tokens are used for low-confident positions. The decision rule can be expressed in Equation (6) below.

$$y_{i-1} = \begin{cases} y_{i-1} & \text{if } p(y_i | X, y_{<i}) \leq 0.9 \\ \hat{y}_{i-1} & \text{if } 0.9 < p(y_i | X, y_{<i}) \leq 0.95 \\ y_{rand} & \text{if } p(y_i | X, y_{<i}) > 0.95 \end{cases} \quad (6)$$

When the probability of the ground truth token at step i in the first decoder is no greater than 0.9, the ground truth token y_{i-1} is chosen as input for the second decoder to reinforce the teacher forcing. When the probability is between 0.9 and 0.95, the token with the maximum probability at step $i-1$ is used to simulate the model prediction in inference. When the probability is larger than 0.95, a token randomly sampled from the target sentence is used.

The output from the second decoder is used with contrastive learning for the Obj_{Rec} and Obj_{CRec} .

To meet the Obj_{Rec} from Inequality (2), we use the function below to formulate the *recovery loss*:

$$\mathcal{L}_{Rec} = \max\{0, \rho + \log p(\hat{y}_i | X, y_{<i-1}, \hat{y}_{i-1}) - \log p(y_i | X, y_{<i-1}, \hat{y}_{i-1})\}. \quad (7)$$

We use the function below to formulate the loss for the Obj_{CRec} (Inequality 3) to *constrain recovery*:

$$\mathcal{L}_{CRec} = \max\{0, \rho + \log p(y_i | X, y_{<i-1}, \hat{y}_{i-1}) - \log p(y_i | X, y_{<i})\}. \quad (8)$$

The overall loss function is a weighted sum of three components:

$$\mathcal{L} = \frac{\mathcal{L}_{MLE} + \alpha \mathcal{L}_{Rec} + \alpha \mathcal{L}_{CRec}}{1 + 2\alpha}, \quad (9)$$

where α is a hyperparameter as the weight.

4 Experiments

4.1 Datasets

Our experiments use the corpora from WMT¹. Wang and Sennrich (2020) claim that the methods reducing exposure bias with sequence-level objectives such as MRT can particularly enhance the model’s resilience to domain shift. To evaluate this claim, we conduct Out-Of-Domain (OOD) tests on De–En and Ru–En.

For De–En, we use Europarl v7, News-commentary-v12, and Common Crawl for training, Newstest2014 for validation, and Newstest2021 and EMEA² for in-domain and OOD testing respectively.

For Ru–En and En–Ru, we use ParaCrawl v9, News-commentary-v10, and Common Crawl for training, Newstest2014 for validation, and Newstest2021 for in-domain testing. The OOD tests for

¹<http://www.statmt.org>

²<https://opus.nlpl.eu/EMEA.php>

Ru–En use the test set for the Biomedical Translation Task in WMT22³.

These original datasets are filtered to remove low-quality data. 350 million sentences are randomly selected with the conditions below:

- The length of source and target sentences are within the range of 5 to 300.
- The disparity between the source and target sentence length does not exceed five times.

The number of sentence pairs in the final training sets for each language pair is: De–En 2.6 million, Ru–En 2.9 million, En–Ru 2.9 million.

4.2 Models

We compare our method to the vanilla Transformer model and reimplement five methods aiming at mitigating exposure bias for comparison.

- *TX* is the vanilla Transformer.
- *SS* (Mihaylova and Martins, 2019) is a typical Scheduled Sampling method based on 2-pass decoding with Transformer. We use Inverse Sigmoid Decay for scheduling since it performs better than other scheduling algorithms according to Liu et al. (2021).
- *CASS* (Liu et al., 2021) is Confidence-Aware Scheduled Sampling using the best configuration in their paper, which outperforms *TFN*, *MIXER*, and *MRT* in their experiments.
- *TFN* (Goodman et al., 2020) uses 2 stacking decoders and combine their loss functions. According to their paper’s recommendation, we use 0.4 as the second decoder’s weight and shared parameter for both decoders.
- *MIXER* (Ranzato et al., 2016): Our implementation follows Kiegl and Kreutzer (2021).
- *MRT* (Shen et al., 2016): We use 4 candidates and do not include the gold reference, same as Wang and Sennrich (2020).

Our method is denoted as *TCL* (*Token-level Contrastive Learning*). The margin ρ for the contrastive learning is set to 0.01. This means that the probability of a negative sample is allowed to reach

³<https://www.statmt.org/wmt22/biomedical-translation-task.html>

99% of the probability of a positive sample maximally. We conducted preliminary experiments on the weight α in the loss function. The models with $\alpha = 0.5$ got bad performance. Our results in this paper are from experiments using models trained with $\alpha = 0.1$.

Our implementation is based on the Fairseq toolkit (Ott et al., 2019) with a typical configuration⁴ similar to the original Transformer (Vaswani et al., 2017). We use the BPE (Sennrich et al., 2016) mode in SentencePiece⁵ for subwords with 32,000 updates and use a shared vocabulary for source and target. We use beam search for decoding. The beam size is 4.

Our experiments are based on Transformer Base (about 60 million parameters). Both the dropout rate and Label Smoothing are set to 0.1 for all models.

The models for vanilla Transformer *TX* are trained for a minimum of 20 epochs, stopping if the validation loss does not decrease for 20 consecutive epochs. The other baseline methods and *TCL* use these vanilla Transformer models as pre-trained models for finetuning. During finetuning, we adopt the same early stop policy as Choshen et al. (2019), where the process is terminated if the validation loss does not decrease for ten consecutive epochs.

The token-level methods (*CASS*, *CASS*, *TFN*, and *TCL*) have similar speeds in training. It takes about 30 minutes to finish one epoch with 10 GPUs. The sequence-level methods (*Mixer* and *MRT*) are much slower since they use online samples during training. It takes *MIXER* and *MRT* about 10 hours and 14 hours to finish one epoch with 10 GPUs, respectively. This result is consistent with the experiments in Edunov et al. (2018). They find that online sampling methods can be 26 times slower than the corresponding offline methods.

We do significance tests for token-level methods. We train models with five different seeds (1–5) and report the mean and standard error over these independent runs. We use the default seed (1) in Fairseq for other experiments. We do not have significance tests for the sequence-level methods (*MIXER* and *MRT*) since they are too slow.

All GPUs that we use are Nvidia GF1080Ti.

⁴https://github.com/facebookresearch/fairseq/tree/main/examples/scaling_nmt

⁵<https://github.com/google/sentencepiece>

Metrics	De-En			Ru-En			En-Ru		
	BLEU	Meteor	Comet	BLEU	Meteor	Comet	BLEU	Meteor	Comet
TX	27.57	49.72	75.01	30.15	49.43	74.93	15.87	29.13	63.97
SS	27.78±.08	49.76±.12	75.16±.01	30.44±.11	49.64±.13	75.16±.07	16.78±.11	30.54±.24	65.95±.34
CASS	27.86±.18	49.74±.07	75.26±.06	30.59±.16	49.85±.10	75.39±.02	17.10±.28	31.08±.05	66.36±.49
TFN	27.62±.23	49.63±.19	75.16±.09	30.44±.10	49.74±.07	75.33±.09	17.04±.18	30.87±.30	66.62±.09
MIXER	27.84	49.74	75.33	30.03	49.67	75.36	17.65	31.64	66.77
MRT	27.41	49.52	75.29	30.39	49.69	75.07	17.15	31.29	66.04
TCL	28.10±.16	49.94±.13	75.33±.07	30.59±.17	49.81±.13	75.50±.18	17.35±.16	31.56±.24	66.83±.13
Δ (-TX)	0.53	0.22	0.32	0.44	0.38	0.57	1.48	2.43	2.86

Table 1: Performance of different methods for the in-domain tests (Newstest2021). We report mean and standard error over five independent training runs with seeds 1–5 for the token-level methods. The scores of TCL and those better than TCL are highlighted in **Bold**. Δ is the gain of TCL compared to TX.

Metrics	De-En			Ru-En		
	BLEU	Meteor	Comet	BLEU	Meteor	Comet
TX	25.75	41.62	67.93	34.94	52.01	74.91
SS	26.17±.14	42.09±.07	68.13±.08	35.66±.06	52.51±.17	75.20±.10
CASS	26.32±.12	42.03±.07	68.23±.09	35.54±.15	52.39±.25	75.28±.12
TFN	26.41±.08	42.04±.06	68.32±.07	35.85±.08	52.57±.13	75.23±.09
MIXER	26.62	42.20	68.50	35.66	52.22	75.18
MRT	26.36	42.05	68.15	35.39	52.55	75.22
TCL	26.62±.20	42.17±.19	68.34±.07	35.82±.11	52.61±.07	75.24±.07
Δ (-TX)	0.87	0.55	0.41	0.88	0.60	0.33

Table 2: Performance of different methods for out-of-domain (OOD) tests. Denotations are the same as Table 1.

4.3 Evaluation and Results

We use BLEU, Meteor, and Comet for evaluation. For BLEU, We use SacreBLEU⁶ (Post, 2018)⁷. For Meteor⁸, we use version 1.5. For Comet⁹, we use the *wmt22-comet-da* model, which scales the scores between 0 and 1. Scores for all metrics are multiplied by 100.

Table 1 and Table 2 illustrate the performance of methods for in-domain test sets (Newstest2021) and out-of-domain test sets, respectively.

TCL outperforms the vanilla Transformer in all tests. TCL gets the best performance among token-level methods in tests except for three cases, highlighted in the tables in **Bold**. The differences in scores between TCL and these three exceptions are very small (less than 0.1).

TCL is ten times more efficient in training compared to those two sequence-level methods. TCL still outperforms those methods in the majority of the tests.

TCL gets larger gains in the OOD tests than in the in-domain tests. This is consistent with the conclusion in Wang and Sennrich (2020). They claim that exposure bias is more influential in domain shift, although their experiment uses the method *MRT*.

Our analysis in Section 5.2 demonstrates that TCL achieves both *recovery* and *constraining recovery* and mitigates the exposure bias. The analysis in Section 5.3 shows the effectiveness of this method by tracking the values of three components in the loss function in training.

5 Analysis

Besides the overall performance, we investigate how these three objectives are met and whether the loss function effectively coordinates these objectives.

We start by using the *prefix switching* method. Then, we directly measure how often the three objectives in Subsection 3.1 are *NOT* met during decoding for each method. Finally, we verify the effectiveness of the loss function in Subsection 3.3 by monitoring how the values of these components

⁶<https://github.com/mjpost/sacreBLEU>

⁷case.mixed+numrefs.1+smooth.exp+tok.l3a+version.2.3.1

⁸<http://www.cs.cmu.edu/~alavie/METEOR/>

⁹<https://github.com/Unbabel/COMET>

	De-En		Ru-En		En-Ru	
	Prefix	Normal	Prefix	Normal	Prefix	Normal
TX	41.37	27.57	42.87	30.15	30.25	15.87
SS	41.20	27.75	43.28	30.20	30.91	16.86
CASS	41.16	27.70	43.42	30.35	31.01	17.19
TFN	41.77	27.25	43.47	30.30	30.83	17.29
MIXER	41.40	27.84	43.43	30.03	30.54	17.65
MRT	40.96	27.41	43.42	30.39	30.83	17.15
TCL	41.65	28.48	43.44	30.39	30.79	17.33

Table 3: The inconsistency between *prefix switching* test (denoted as *Prefix*) and normal tests. Best BLEU scores are highlighted in **Bold**.

in our loss function change during training TCL and its variants.

5.1 Using Prefix Switching to Quantify Exposure Bias Is Not Reliable

Prefix Switching is often used to quantify exposure bias (Wu et al., 2018; Korakakis and Vlachos, 2022). We use various lengths of ground truth tokens as prefixes and measure the average quality of the part of the sequence from the model’s prediction. The length of the prefix varies from 1 to $N-1$, where N is the length of the reference. After decoding, we measure the average *sentence-BLEU* scores of the prediction part of sequences. If the length of a prediction part is shorter than 4, it is not considered for the average.

Table 3 shows the results for three language pairs on the in-domain test sets using the Prefix Switching and the normal tests. In the normal tests, there are no ground-truth prefixes during decoding.

The results of these two tests are inconsistent. For example, TFN gets the best BLEU score in De-En in *Prefix Switching* testing. But it gets a score lower than the vanilla Transformer in the normal test. It reflects that using prefix switching to quantify the exposure bias may not be reliable. This issue requires further investigation.

5.2 Analysis if Three Objectives Are Met or Not

We directly detect how many times these three objectives (Obj_{MLE} , Obj_{Rec} and Obj_{CRec}) in Subsection 3.1 are met or not in decoding.

Similar to prefix switching, we use various lengths of ground truth as prefixes to the decoder. In this experiment, we only need to monitor one or two steps of decoding, not requiring the decoder to finish a prediction with an *End-of-Sentence (EOS)*.

Assume that N is the length of the gold reference in subwords and k is the length of the prefix which enumerates between 0 and N .

We need a set of *non ground-truth* tokens, \hat{y}_{k-1} , to test Obj_{Rec} and Obj_{CRec} . \hat{y}_{k-1} is corresponding to \hat{y}_{i-1} in Inequality (2) and (3). It is intractable to enumerate all tokens in the vocabulary. We choose the m tokens with the top probabilities from outputs at the step $k-2$ and test with each of them by appending it to the ground truth prefixes $y_{<k-1}$ for decoding at step $k-1$. The ground truth token is taken out if it is in this top- m token set. We use $m \in \{1, 5, 10\}$ since the size of the beam search is usually not greater than 10 in practice.

Once \hat{y}_{k-1} is selected, the decoder uses $(X, y_{<k-1}, y_{k-1})$ and $(X, y_{<k-1}, \hat{y}_{k-1})$ as inputs *respectively* and get both $p(y_k|X, y_{<k})$ and $p(y_k|X, y_{<k-1}, \hat{y}_{k-1})$. If $p(y_k|X, y_{<k})$ is the maximum in its decoding step, we can tell that Obj_{MLE} is met. If $p(y_k|X, y_{<k-1}, \hat{y}_{k-1})$ is the maximum in its step, Obj_{Rec} is met.

We use the joint probability of *bi-gram* to test Obj_{CRec} , the third (missing) objective, since the total probability of the sequence is used in decoding. Inequality (10) below is the criterion:

$$\begin{aligned}
 & p(y_k, y_{k-1}|X, y_{<k-1}) = \\
 & p(y_k|X, y_{<k}) * p(y_{k-1}|X, y_{<k-1}) \\
 & > \\
 & p(y_k, \hat{y}_{k-1}|X, y_{<k-1}) = \\
 & p(y_k|X, y_{<k-1}, \hat{y}_{k-1}) * p(\hat{y}_{k-1}|X, y_{<k-1})
 \end{aligned} \tag{10}$$

Table 4, 5, and 6 illustrates results for different methods for De-En, Ru-En and En-Ru respectively. The event of *Not Met* is counted for each step for each objective. When the number of *non ground-truth* tokens (m) is larger than 1, such an event may happen more than once at one step. The

	Obj_{MLE}	Obj_{Rec}			Obj_{CRec}		
		Top1	Top5	Top10	Top1	Top5	Top10
TX	0.316	0.277	3.087	6.852	0.125	0.362	0.534
SS	0.314	0.275	3.091	6.855	0.127	0.362	0.530
CASS	0.314	0.275	3.078	6.806	0.127	0.366	0.538
TFN	0.313	0.275	3.117	6.933	0.138	0.388	0.566
TCL	0.314	0.275	3.088	6.868	0.123	0.354	0.521
MIXER	0.317	0.279	3.101	6.887	0.125	0.363	0.532
MRT	0.314	0.274	3.080	6.842	0.127	0.353	0.513

Table 4: Failure rates of three objectives for De–En. **Smaller is better.** The smallest ones are highlighted in **Bold**. The values in this table are how often the objective is *NOT* met, divided by the total number of tests (24760 in this case). $Top-m$ denotes that the number of *non ground-truth* tokens (\hat{y}_{k-1}) used in test is m . CASS has a larger failure rate for the third objective Obj_{CRec} than the vanilla Transformer. This result reflects that CASS has enhanced recovery *too much* that it deviates from the ground truth. TCL is the only token-level method with lower failure rates for all objectives than the vanilla Transformer. The two sequence-level methods are not supposed to have the deviation issue, but they are tested here for reference.

	Obj_{MLE}	Obj_{Rec}			Obj_{CRec}		
		Top1	Top5	Top10	Top1	Top5	Top10
TX	0.288	0.254	3.078	6.891	0.104	0.281	0.394
SS	0.285	0.251	3.083	6.891	0.104	0.276	0.388
CASS	0.286	0.250	3.068	6.838	0.106	0.285	0.400
TFN	0.284	0.249	3.103	6.955	0.115	0.303	0.418
TCL	0.285	0.251	3.067	6.865	0.103	0.275	0.387
MIXER	0.285	0.252	3.078	6.892	0.104	0.275	0.388
MRT	0.285	0.251	3.079	6.874	0.102	0.274	0.380

Table 5: Failure rates of three objectives for Ru–En. **Smaller is better.** The denotations are the same as Table 4. The total number of tests is 27828 in this case.

	Obj_{MLE}	Obj_{Rec}			Obj_{CRec}		
		Top1	Top5	Top10	Top1	Top5	Top10
TX	0.379	0.356	3.464	7.559	0.149	0.492	0.776
SS	0.375	0.351	3.462	7.560	0.145	0.487	0.774
CASS	0.373	0.348	3.431	7.492	0.151	0.504	0.794
TFN	0.373	0.353	3.464	7.558	0.159	0.519	0.813
TCL	0.376	0.352	3.451	7.536	0.144	0.486	0.770
MIXER	0.375	0.352	3.467	7.577	0.140	0.467	0.740
MRT	0.373	0.349	3.448	7.540	0.146	0.479	0.757

Table 6: Failure rates of three objectives for En–Ru. **Smaller is better.** The denotations are the same as Table 4. The total number of tests is 42442 in this case.

total number of events is then divided by the number of steps (for example, 24760 in De–En). The results are the *average failure rate* per token.

These tables show that CASS has the lowest failure rates for the second objective Obj_{Rec} in both De–En and En–Ru. CASS also gets relatively low failure rates for this objective in Ru–En. These results demonstrate that CASS successfully enhances the recovery capability. However, CASS has larger failure rates for the third objective Obj_{CRec} than the vanilla Transformer in all three

language pairs. This result reveals that CASS has enhanced recovery *too much* that it deviates from the ground truth, which is the side effect described in Subsection 3.1.

Our method TCL gets the lowest failure rate for the third objective Obj_{CRec} among the token-level methods in all tests. Furthermore, TCL is the only token-level method with lower failure rates for all objectives than the vanilla Transformer in Ru–En and En–Ru. It achieves a *pareto optimality* in the sense of improvement on both objectives: *recovery*

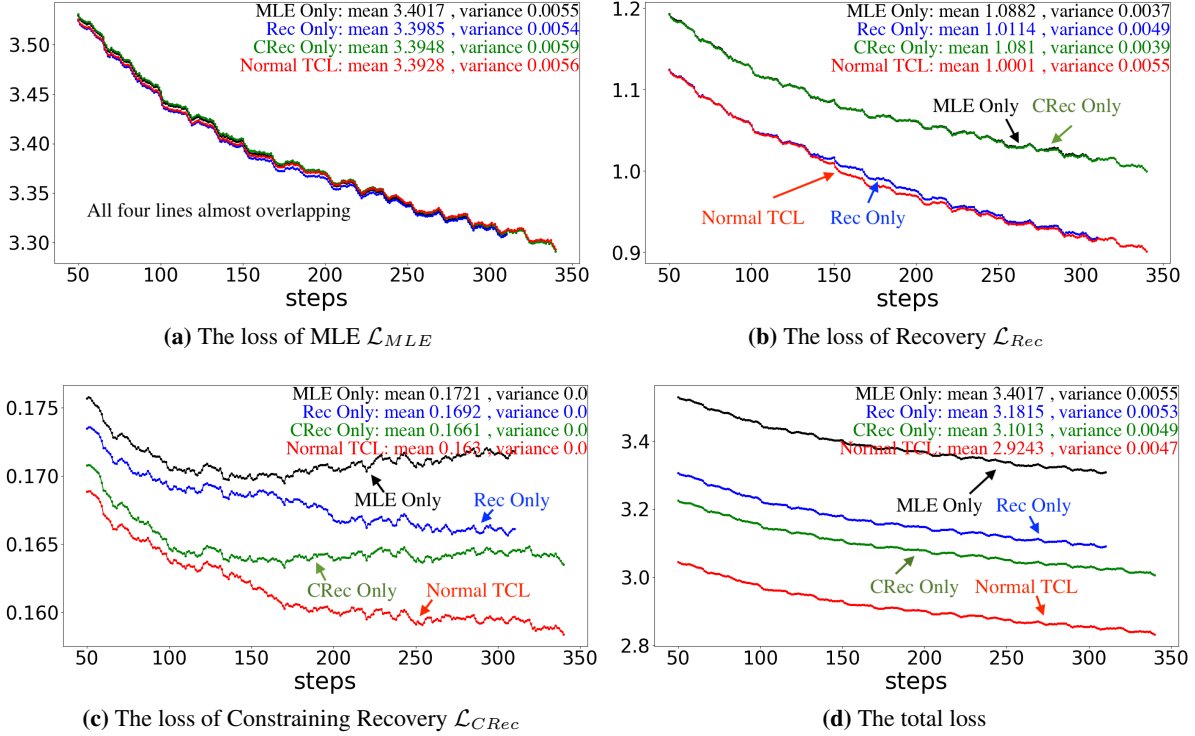


Figure 2: Investigate the values of components in TCL’s loss function for De–En in training. *Rec Only* denotes the model trained without applying the loss component \mathcal{L}_{CRec} . *CRec Only* denotes the model trained without applying the loss component \mathcal{L}_{Rec} . *MLE Only* denotes the model trained without applying both \mathcal{L}_{Rec} and \mathcal{L}_{CRec} .

and *constraining recovery*. These results demonstrate that the exposure bias is mitigated by our method.

The sequence-level methods do not have the deviation issue discussed in this paper since they use sequence-level objectives in training. Their results are included in these failure rate tests for reference only. The results show that they perform well in this test, reflecting their effectiveness in mitigating exposure bias, although these methods are much slower than the token-level methods.

5.3 Effectiveness of Loss Components

There are three components in our loss function in Equation (9): \mathcal{L}_{MLE} , \mathcal{L}_{Rec} , and \mathcal{L}_{CRec} . We evaluate the effectiveness of these components by tracking their loss values in training TCL and its three variants by turning off one or two components. We use α_1 and α_2 to denote the weights for \mathcal{L}_{Rec} and \mathcal{L}_{CRec} , respectively.

- *Normal TCL*: $\alpha_1 = \alpha_2 = 0.1$
- *Rec Only* (recovery): $\alpha_1 = 0.1, \alpha_2 = 0$
- *CRec Only* (constraining recovery): $\alpha_1 = 0, \alpha_2 = 0.1$

- *MLE Only*: $\alpha_1 = \alpha_2 = 0$

Figure 2 illustrates how each loss component’s values vary in training for De–En. These values are reported every 100 updates during training and smoothed by taking the average with their ten right and left neighbors.

Figure 2a shows that the values of \mathcal{L}_{MLE} for four models are almost the same. This component is not influenced by other two components.

Figure 2b shows the recovery loss \mathcal{L}_{Rec} . Even for the model *MLE Only* without \mathcal{L}_{Rec} and \mathcal{L}_{CRec} , this loss decreases in training. This implies that models increase self-recovery capability during training even if no extra means are used to enhance it. This result supports the conclusion from He et al. (2021), although enhancing the recovery capability may not be enough to deny exposure bias’s negative impact. The blue and red lines (*Rec Only* and *Normal TCL*) with the recovery component get smaller values than the other two models without this component. This illustrates that this component in the loss function effectively increases the capability of recovery.

Figure 2c shows the values of \mathcal{L}_{CRec} (constraining recovery). Similar to the values of \mathcal{L}_{Rec} in Figure 2b, even for the model *MLE Only* without \mathcal{L}_{Rec}

and \mathcal{L}_{CRec} , this loss decreases in training. The green and red lines (*CRec Only* and *Normal TCL*) with the component \mathcal{L}_{CRec} get smaller values than the other two models without this component. This implies that using this component in the loss function effectively reduces the \mathcal{L}_{CRec} .

This loss surprisingly increases after a period of decreasing in training for *MLE Only* and *CRec Only*. This is the consequence of increasing the capability of self-recovery shown in Figure 2b with or without \mathcal{L}_{Rec} . The increasing of $p(y_i|X, y_{<i-1}, \hat{y}_{i-1})$ may result in the increase of values of \mathcal{L}_{CRec} according to its definition in Equation (8). Current token-level methods that maximizes $p(y_i|X, y_{<i-1}, \hat{y}_{i-1})$ may make this contradiction more severe.

Figure 2d shows the total loss.

Table 7 shows the ablation tests using the BLEU scores for *Rec Only* (recovery) and *CRec Only* (constraining recovery) models compared to the vanilla Transformer and the normal TCL models. *Rec Only* gets worse performance than the vanilla Transformer. *CRec Only* have some gains. The normal TCL that combines these components gets extra improvement. Table 8 in Appendix A illustrates the results for En–Ru, and they lead to the same conclusion.

Metrics	De–En		
	BLEU	Meteor	Comet
Vanilla Transformer (TX)	27.57	49.72	75.01
Rec Only	27.23	49.29	75.27
Δ (-TX)	-0.34	-0.43	0.26
CRec Only	27.82	49.85	75.40
Δ (-TX)	0.25	0.13	0.39
TCL	28.48	50.20	75.55
Δ (-TX)	0.91	0.48	0.54

Table 7: Ablation tests. *Rec Only* (recovery) and *CRec Only* (constraining recovery) models compared to the vanilla Transformer and normal TCL models.

6 Conclusion

Current token-level methods addressing exposure bias may have a side effect: A sequence with errors may have a larger probability than the ground truth. Consequently, the generated sequence may deviate from the ground truth. Our experiments verify this side effect. We discover a missing objective behind this side effect that can explicitly constrain the recovery in a scope that does not im-

pact the ground truth. We propose token-level contrastive learning to coordinate three objectives in the loss function: the original MLE, recovery from errors, and constraining the recovery in a scope not to exceed the ground truth. Experimental results on three language pairs show that our method outperforms the vanilla Transformer and five methods aiming at mitigating exposure bias. Empirical analysis demonstrates that this method achieves a Pareto optimality compared with the vanilla Transformer. It is also verified that each component in our loss function effectively improves the model in training.

References

- Arora, Kushal, Layla El Asri, Hareesh Bahuleyan, and Jackie Chi Kit Cheung. 2022. Why exposure bias matters: An imitation learning perspective of error accumulation in language generation. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 700–710.
- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Bengio, Samy, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. 2015. Scheduled sampling for sequence prediction with recurrent neural networks. *Advances in neural information processing systems*, 28.
- Chiang, Ting-Rui and Yun-Nung Chen. 2021. Relating neural text degeneration to exposure bias. In *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 228–239.
- Choshen, Leshem, Lior Fox, Zohar Aizenbud, and Omri Abend. 2019. On the weaknesses of reinforcement learning for neural machine translation. In *International Conference on Learning Representations*.
- Edunov, Sergey, Myle Ott, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. 2018. Classical structured prediction losses for sequence to sequence learning. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 355–364.
- French, Robert M. 1999. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3(4).
- Goodman, Sebastian, Nan Ding, and Radu Soricut. 2020. Teaform: Teacher-forcing with n-grams. In *Proceedings of the 2020 Conference on Empirical*

- Methods in Natural Language Processing (EMNLP)*, pages 8704–8717.
- He, Tianxing, Jingzhao Zhang, Zhiming Zhou, and James Glass. 2021. Exposure bias versus self-recovery: Are distortions really incremental for autoregressive text generation? In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5087–5102.
- Huszár, Ferenc. 2015. How (not) to train your generative model: Scheduled sampling, likelihood, adversary? *arXiv preprint arXiv:1511.05101*.
- Kiegeland, Samuel and Julia Kreutzer. 2021. Revisiting the weaknesses of reinforcement learning for neural machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1673–1681.
- Korakakis, Michalis and Andreas Vlachos. 2022. Improving scheduled sampling with elastic weight consolidation for neural machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7247–7258.
- Liu, Yijin, Fandong Meng, Yufeng Chen, Jinan Xu, and Jie Zhou. 2021. Confidence-aware scheduled sampling for neural machine translation. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2327–2337.
- Liu, Yixin, Pengfei Liu, Dragomir Radev, and Graham Neubig. 2022. Brio: Bringing order to abstractive summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2890–2903.
- Mihaylova, Tsvetomila and André FT Martins. 2019. Scheduled sampling for transformers. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, pages 351–356.
- Ott, Myle, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53.
- Pan, Xiao, Mingxuan Wang, Liwei Wu, and Lei Li. 2021. Contrastive learning for many-to-many multilingual neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 244–258.
- Post, Matt. 2018. A call for clarity in reporting bleu scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191.
- Ranzato, Marc’Aurelio, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2016. Sequence level training with recurrent neural networks. In *4th International Conference on Learning Representations, ICLR 2016*.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725.
- Shen, Shiqi, Yong Cheng, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. 2016. Minimum risk training for neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1683–1692.
- Su, Yixuan, Tian Lan, Yan Wang, Dani Yogatama, Lingpeng Kong, and Nigel Collier. 2022. A contrastive framework for neural text generation. *Advances in Neural Information Processing Systems*, 35:21548–21561.
- Sun, Shichao and Wenjie Li. 2021. Alleviating exposure bias via contrastive learning for abstractive text summarization. *arXiv preprint arXiv:2108.11846*.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Wang, Chaojun and Rico Sennrich. 2020. On exposure bias, hallucination and domain shift in neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3544–3552.
- Wu, Lijun, Xu Tan, Di He, Fei Tian, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. 2018. Beyond error propagation in neural machine translation: Characteristics of language also matter. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3602–3611.
- Yang, Zonghan, Yong Cheng, Yang Liu, and Maosong Sun. 2019. Reducing word omission errors in neural machine translation: A contrastive learning approach. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6191–6196.
- Zhang, Wen, Yang Feng, Fandong Meng, Di You, and Qun Liu. 2019. Bridging the gap between training and inference for neural machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4334–4343.

A Ablation tests for En–Ru

Table 8 shows the ablation tests for En–Ru. Both *Rec Only* and *CRec Only* have some gains. The normal TCL that combines these components gets extra improvement.

Metrics	En–Ru		
	BLEU	Meteor	Comet
Vanilla Transformer (TX)	15.87	29.13	63.97
Rec Only	16.33	29.71	65.05
Δ (-TX)	0.46	0.58	1.08
CRec Only	16.65	30.86	65.92
Δ (-TX)	0.78	1.73	1.95
TCL	17.33	31.77	67.02
Δ (-TX)	1.46	2.64	3.05

Table 8: Ablation tests. *Rec Only* (recovery) and *CRec Only* (constraining recovery) models compared to the vanilla Transformer and normal TCL models.

Chasing COMET: Leveraging Minimum Bayes Risk Decoding for Self-Improving Machine Translation

Kamil Guttman^{*1}, Mikołaj Pokrywka^{*1}, Adrian Charkiewicz¹, Artur Nowakowski^{1,2}

¹ Lanigo, Poznań, Poland

² Faculty of Mathematics and Computer Science, Adam Mickiewicz University, Poznań, Poland
{name}. {surname}@lanigo.com

Abstract

This paper explores Minimum Bayes Risk (MBR) decoding for self-improvement in machine translation (MT), particularly for domain adaptation and low-resource languages. We implement the self-improvement process by fine-tuning the model on its MBR-decoded forward translations. By employing COMET as the MBR utility metric, we aim to achieve the reranking of translations that better aligns with human preferences. The paper explores the iterative application of this approach and the potential need for language-specific MBR utility metrics. The results demonstrate significant enhancements in translation quality for all examined language pairs, including successful application to domain-adapted models and generalisation to low-resource settings. This highlights the potential of COMET-guided MBR for efficient MT self-improvement in various scenarios.

1 Introduction

Machine translation (MT) bridges the gap between languages, fostering global communication and information exchange. However, achieving high-quality translations across diverse languages and domains remains a significant challenge, especially for low-resource languages where limited training data hinders model performance. Even in well-resourced settings, continuous improvement

and adaptation to specific domains are ongoing research efforts.

This paper explores the potential of Minimum Bayes Risk (MBR) decoding (Kumar and Byrne, 2004) as a self-improvement strategy for MT models. MBR decoding leverages the model’s predictions to select the best translation from a set of candidates, potentially improving overall translation quality.

We employ COMET (Rei et al., 2020) as the utility function in MBR decoding and rerank candidate translations generated by an MT model. This approach creates a synthetic parallel dataset from monolingual data in the source language, enabling further model self-improvement.

This study examines the effectiveness of MBR decoding for self-improvement in three language pairs: English–German (high-resource), Czech–Ukrainian (low-resource), and English–Hausa (low-resource). For English–German, the focus is on the biomedical domain, incorporating additional monolingual data, while for Czech–Ukrainian, self-improvement is explored using only the training data translated by the model and reranked through MBR decoding. We further investigate the potential of iterative self-improvement with MBR decoding in both English–German and Czech–Ukrainian language pairs. Finally, in the case of English–Hausa, we compare the use of COMET, a massively multilingual metric, with a metric specifically tailored to African languages i.e. AfriCOMET (Wang et al., 2023).

To determine the optimal configuration for MBR decoding, we investigate two decoding algorithms and various numbers of translation candidates.

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

*Equal contribution

2 Related Work

MBR and QE reranking with neural metrics

MBR decoding, a technique commonly used in Statistical Machine Translation (SMT), has gained traction in Neural Machine Translation (NMT) in recent years. Freitag et al. (2022) proposed using reference-based metrics, such as BLEURT (Selam et al., 2020a) and Quality Estimation (QE) models, such as COMET-QE (Rei et al., 2021) for reranking the set of hypotheses produced by the NMT model.

Similar work by Fernandes et al. (2022) proposed *quality-aware decoding*. They explored various reranking strategies, including the well-performing pre-ranking of the set of hypotheses with QE models before passing them into MBR decoding. They found that using MERT-tuned (Och, 2003) reranker, where multiple QE metrics and model log-likelihood scores are linearly combined with learned weights to maximize a reference-based metric on a validation set shows improvements over the baseline.

Amrhein and Sennrich (2022) used MBR decoding to identify biases and weaknesses in COMET, where they found that the early COMET models are not sufficiently sensitive to discrepancies in numbers and named entities.

MBR decoding performance is heavily dependent on the number of samples and the sampling strategy. Freitag et al. (2023) investigated various sampling strategies and found that epsilon sampling outperformed others. This sampling method discards tokens with a probability below a certain threshold (epsilon), guaranteeing that each token in the final sample has a fair chance of being included. The approach is particularly effective when generating a large set of samples, as it inherently yields greater sample diversity compared to beam search.

Vernikos and Popescu-Belis (2024) introduced QE-fusion, a method that combines spans from different candidates sampled from a model using QE metrics. They found that the method consistently improves translation quality in terms of neural evaluation metrics, especially if applied to LLM due to their ability to generate diverse outputs.

Due to its ease of implementation and use, MBR and QE reranking have been successfully applied in machine translation shared tasks, as demonstrated by the results in several stud-

ies (Nowakowski et al., 2022; Kudo et al., 2023; Jon et al., 2023). This highlights its potential to significantly improve translation quality.

Model self-improvement Recent research has shown a growing interest in leveraging model outputs for self-improvement. This approach holds significant promise in the case of machine translation, especially for low-resource and domain-specific translation scenarios, where there is access to the source-language data, but the corresponding target-language data is severely limited.

Gulcehre et al. (2023) describes reinforcement self-training (*ReST*) method for language modeling. The method is based on producing a dataset for fine-tuning by sampling from the model (LLM). The samples are then scored with a QE metric. Then, offline reinforcement learning algorithms are applied using a reward-weighted loss based on the QE scores. The method can be applied to all generative learning settings, but the authors focus on its application to machine translation, showing that the method increases translation quality.

Concurrent work by Finkelstein et al. (2023) describes self-tuning NMT models on a set of hypotheses reranked using either MBR, QE, or a combination of the two methods. They also experimented with using LLM as the teacher model, finding that it outperforms using a self-teacher and fine-tuning on references.

Our research expands on recent developments in the field by investigating the use of MBR-based fine-tuning in three key areas. Firstly, we examine its applicability in domain-specific translation tasks, specifically focusing on English–German translation in the biomedical domain. Secondly, we investigate its effectiveness for low-resource translation directions, exemplified by the Czech–Ukrainian language pair. This broadens the scope beyond English-centric language pairs, thus contributing to a more comprehensive analysis of MBR performance across less-represented languages in neural evaluation metrics. Finally, we explore the use of neural QE metrics tailored for specific languages, using AfriCOMET (Wang et al., 2023) as an example.

3 Experiment Overview

3.1 Model Self-Improvement

The self-improvement process leverages MBR decoding to guide the model to select high-quality translations according to the utility function. The process consists of 3 steps:

Step 1: Sample Generation Using beam search decoding with beam size equal to N , generate N translation candidates using the base model for each source sentence. While Freitag et al. (2023) suggested that epsilon sampling might yield better results with MBR decoding, it typically requires reranking a significantly larger number of translation candidates, which becomes computationally expensive for processing large datasets. Beam search, on the other hand, allows for generating a smaller set of high-quality candidates while providing sufficient data for effective MBR decoding.

Step 2: MBR Decoding Select a single translation for each source sentence from the list of candidates through MBR decoding utilizing COMET to guide the selection towards high-quality translations. For an efficient implementation of the MBR decoding algorithm, we use the code¹ from the Marian (Junczys-Dowmunt et al., 2018) framework.

Step 3: Model Fine-tuning Fine-tune the base model on the synthetically created dataset. Use COMET as an early stopping metric during training to ensure fitting to this metric.

3.2 English–German

The English–German experiment simulates a real-world domain adaptation scenario. In such settings, while a large general-purpose parallel corpus might be available, the specific domain often lacks extensive parallel data. To address this challenge, we leveraged both a smaller parallel dataset and a larger monolingual dataset in the source language containing biomedical terminology.

To leverage the monolingual data in the source language we propose a two-step approach:

1. Fine-Tuning: We fine-tune a general-purpose English–German model on a small parallel biomedical dataset.

¹<https://github.com/marian-nmt/marian-dev/tree/master/scripts/mbr>

2. Self-improvement: To enhance the model performance in the biomedical domain, we incorporate a larger monolingual biomedical dataset during the self-improvement process. This involves creating a synthetic parallel dataset via MBR decoding and subsequently fine-tuning the biomedical translation model on the generated data.

To assess the robustness of the self-improvement method, we conducted an additional experiment in which we applied this method to a model that was fine-tuned to the biomedical domain using general domain data for MBR decoding. This evaluated whether the model would retain its translation capabilities in the biomedical domain despite improvements based solely on out-of-domain data.

3.3 Czech–Ukrainian

The Czech–Ukrainian experiment addresses the challenge of machine translation between two low-resource languages. We aimed to evaluate whether self-improvement through MBR decoding leads to an increase in the overall translation quality when applied to language pairs that do not involve English, which typically dominate machine translation research.

In this setting, we used only the parallel data set without incorporating any additional monolingual data. To employ MBR decoding in this data-scarce environment, we directly translated the entire source side of the parallel dataset using the baseline translation model. This created a set of synthetic candidate translations, which were then reranked through MBR decoding.

In contrast to our English–German experiments where we incorporated external monolingual data, this setup explored self-improvement without relying on additional datasets. We achieved this by solely leveraging the information present within the data of the base model. This demonstrates the potential for self-improvement even in resource-constrained scenarios.

3.4 English–Hausa

The English–Hausa experiment delves into the critical question of how the choice of a quality evaluation metric influences the effectiveness of self-improvement with MBR decoding. We explored the impact of language coverage in the evaluation metric by comparing two approaches:

- MBR decoding with WMT22 COMET: utilizing the *wmt22-comet-da* model, which has been trained on direct assessments between a diverse set of language pairs.
- MBR decoding with AfriCOMET: using AfriCOMET-STL, a novel COMET-like metric specifically designed for evaluating translations to and from multiple African languages, including Hausa.

The objective of this study was to investigate the effect of language contribution in the neural evaluation metric on the quality of translations decoded using MBR. The comparison of these two approaches specifically addresses whether self-improvement guided by the WMT22 COMET metric, which is trained on a diverse range of language pairs, can effectively generalize to low-resource language pairs. Furthermore, we explore the potential need to use language-specific metrics, such as AfriCOMET-STL for Hausa, to achieve better performance in such scenarios.

3.5 Iterative MBR Self-Improvement

Following the initial self-improvement through MBR decoding, we explored the possibility of applying it iteratively to further enhance the model’s translation quality.

We started each iteration by selecting the best model checkpoint based on the WMT22 COMET metric on the validation set. Next, we performed MBR decoding on the entire training set using this checkpoint, generating a new iteration of the synthetic training set. Finally, we resumed the training of the model using the new training set, starting from the previously selected checkpoint.

The iterative process was repeated until a decrease was observed in the evaluation scores of metrics other than WMT22 COMET. In the case of English–German biomedical translation, the process was continued until the model’s quality improved solely on an in-domain test set and decreased on a general domain test set, as this could indicate potential overfitting to the biomedical domain.

4 Experimental Setup

4.1 Data Filtering

We filtered the general training data using the following heuristic filters:

- average length of words in each sentence (character-wise) ≤ 15 ;
- number of characters in each sentence ≤ 500 ;
- digits in a sentence (character-wise) $\leq 15\%$;
- number of characters in the longest word ≤ 28 ;
- number of words in sentence ≤ 100 ;
- Levenshtein distance between source and target sentences ≥ 2 ;
- number of characters in each sentence ≥ 5 ;
- probability that each sentence is in the correct language $\geq 10\%$.

To ensure that each sentence is in the correct language, we have used the fastText LID-201 language identification model (Burchell et al., 2023).

The Bicleaner-AI model (Zaragoza-Bernabeu et al., 2022) is also used to filter the English–German dataset. This tool estimates the likelihood that a sentence pair constitutes a mutual translation. A threshold of 50% is established for the Bicleaner score within this language pair. Bicleaner-AI is not utilized for other language pairs due to the unavailability of open-source models for those languages.

4.2 Vocabulary

We employed SentencePiece (Kudo and Richardson, 2018), a subword tokenization library, to train unigram tokenizers for each language pair in our experiments.

For the English–German and English–Hausa setups, we created a joint vocabulary containing 32,000 subword tokens and tied all embeddings during the training of the MT model. In contrast, for Czech–Ukrainian, due to different scripts (Latin and Cyrillic), we created separate vocabularies of 32,000 subword tokens and tied only the target and output layer embeddings.

4.3 Baseline Model Hyperparameters

For all experiments, we trained Transformer (big) (Vaswani et al., 2017) models using the Marian framework. These models were trained on four NVIDIA A100 GPUs, each equipped with 80GB of VRAM.

Hyperparameter Settings:

- learning rate: $2e-4$;

- learning rate warmup: 8000 updates;
- learning rate decay: inverse square root;
- mini-batch size determined automatically to fit GPU memory;
- early stopping after 10 consecutive validations with no improvement in mean word cross-entropy score.

4.4 Evaluation metrics

We use sacreBLEU (Post, 2018) to calculate BLEU² (Papineni et al., 2002) and chrF³ (Popović, 2015).

We acknowledge the potential for overfitting to the WMT22 COMET⁴ metric used for MBR decoding. Therefore, we extended the evaluation to also include CometKiwi⁵ (Rei et al., 2022), UniTE⁶ (Wan et al., 2022), UniTE-DA⁷ (Rei et al., 2023) and BLEURT-20⁸ (Sellam et al., 2020b).

For the English–Hausa experiments, we additionally calculated scores using AfriCOMET-STL (Wang et al., 2023), which was specifically trained to evaluate translations involving certain African languages.

4.5 English to German

To train the baseline model, we used all corpora from the MTData toolkit (version 0.4.0) (Gowda et al., 2021), excluding the validation sets and the test sets from the available datasets. Our filters described in Section 4.1 reduced the dataset from approximately 800 million sentences to 400 million.

In the context of domain adaptation, we employed the following list of domain data:

1. 40 thousand sentences from biomedical-translation-corpora (Neves et al., 2016);
2. 3 million sentences from Ufal medical corpus shared in WMT23 (Kocmi et al., 2023);

²BLEU signature: nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1

³chrF signature: nrefs:1|case:mixed|eff:yes|nc:6|nw:0|space:no|version:2.3.1

⁴<https://huggingface.co/Unbabel/wmt22-comet-da>

⁵<https://huggingface.co/Unbabel/wmt22-cometkiwi-da>

⁶<https://huggingface.co/Unbabel/unite-mup>

⁷<https://huggingface.co/Unbabel/wmt22-unite-da>

⁸<https://storage.googleapis.com/bleurt-oss-21/BLEURT-20.zip>

3. 2 million sentences from EMEA corpus downloaded from OPUS (Tiedemann and Nygaard, 2004).

After deduplication, we were left with 3 million sentences which we split into two datasets. We considered a scenario with 1 million bilingual parallel sentences and approximately 2 million monolingual sentences in the source language. Khresmoi-dev (Dušek et al., 2017) concatenated with FLORES-200 (NLLB Team et al., 2022) was utilized as the validation set during training. We did not apply any filtering to the domain data.

We used the above data to train the following models:

- Baseline (**Baseline**) – model trained only on data from the MTdata toolkit.
- Baseline + mix-tuning (**Mix-tune**) – fine-tuned **Baseline** model on 1 million in-domain bilingual data concatenated with 1 million general-domain data randomly sampled from the **Baseline** training set.
- Baseline + domain MBR (**Base-domain-mbr**) – fine-tuned **Baseline** model on 2 million domain-specific sentences from MBR-decoded forward translations.
- Mix-tuned + domain MBR (**Mix-tune-domain-mbr**) – fine-tuned **Mix-tune** model on 2 million domain-specific sentences from MBR-decoded forward translations.
- Mix-tuned + MBR-iteration2 (**Mix-tune-domain-mbr-iter2**) – fine-tuned **Mix-tune-domain-mbr** on the 2 million domain-specific sentences from MBR-decoded forward translations.
- Mix tuned + general-MBR (**Mix-tune-general-mbr**) – fine-tuned **Mix-tune** model on 2 million sentences sampled from the general-domain corpora from the **Baseline** training set as MBR-decoded forward translations.

When fine-tuning the **Mix-tune** model, we tailor the learning rate setup to meet specific requirements: learn-rate: 1e-7, lr-decay-inv-sqrt: 16000, lr-warmup: 16000. All remaining fine-tuning procedures employ an adjusted learning rate set to 5e-6.

4.6 Czech to Ukrainian

We leveraged all of the Czech–Ukrainian parallel data from the WMT23 MTData recipe, resulting in approximately 8 million sentence pairs after filtering as described in Section 4.1. We did not include any additional monolingual data in this experiment.

We utilized the FLORES-200 dataset for validation during training, while the WMT22 test set served as an additional benchmark.

We trained the baseline model only on the parallel data, using hyperparameters as described in Section 4.3. Next, we translated the source side of the parallel corpus used in training with our baseline model, saving a list of translation candidates. We performed MBR decoding, selecting the best translation of each set of candidate translations, resulting in a synthetic training dataset.

We investigated the following approaches to leverage the MBR-decoded data for model improvement:

- Standard fine-tuning (**MBR-finetuned**) – we fine-tuned the baseline model on the MBR-decoded data, using a learning rate of $5e-6$.
- Fine-tuning with a high learning rate (**MBR-ft-high-lr**) – we fine-tune the baseline model on MBR-decoded data, using a learning rate of $2e-4$.
- Resuming training with MBR-decoded data (**MBR-resumed**) – we switched the training set to the MBR-decoded version and resumed training, restoring the optimizer state and effectively continuing its training with the improved data.

4.7 English to Hausa

To train the models in the English–Hausa direction, we used data from the WMT shared tasks from previous years. Specifically, we used:

1. 7 million sentences from OPUS;
2. 2.4 million data from the WMT23 African MT Shared Task (Kocmi et al., 2023);
3. 150 thousand sentences from ParaCrawl v8.0 (Bañón et al., 2020).

The deduplication process reduced the data size to approximately 9 million sentences. Following the filtering criteria detailed in Section 4.1, a total

of 3.1 million sentences were retained. We used FLORES-200 for validation during training. After training, we evaluated the model on the FLORES-200 and NTREX test sets.

We took similar steps as in the Czech–Ukrainian experiment, training a baseline model with hyperparameters set as described in Section 4.3. We conducted experiments employing MBR decoding, comparing its performance using two distinct metrics as the utility function:

- WMT22 COMET – based on XLM-RoBERTa (Conneau et al., 2020), covering a diverse set of 100 languages,
- AfriCOMET-STL – based on AfroXLM-RoBERTa (Alabi et al., 2022), covering 17 African languages and 3 high-resource languages.

We investigated the impact of the chosen metric for MBR decoding by training two models using the refined translations:

- **MBR-COMET** – training resumed with the training set switched to the WMT22 COMET MBR-decoded version.
- **MBR-AfriCOMET** – training resumed with the training set switched to the AfriCOMET-STL MBR-decoded version.

5 Results

The statistical significance of the evaluation results is assessed using a paired bootstrap resampling test (Koehn, 2004), involving 1000 resampling trials to confirm the statistical significance of the model improvements ($p < 0.05$).

5.1 Number of translation samples and search algorithm

To determine the optimal setup for MBR decoding, we conducted experiments involving the translation and evaluation of chosen test sets with various MBR decoding sample sizes and two decoding algorithms. This approach offers the advantages of being both representative and computationally efficient compared to training MT models on the entire MBR-decoded training set.

We evaluated two decoding algorithms – beam search and top-k. For the top-k setup, we experimented with temperature values of 0.1 and 1, keeping the k parameter equal to 10. These choices

were based on the work done by Freitag et al. (2023). To determine the best number of samples for MBR decoding we conducted experiments with the following numbers of samples: 10, 25, 50, 100, 200, 300, 400, 500.

Firstly we noted that beam search is the preferred option, given its high scores and greater stability across different metric results, as observed in Figure 1 and 2. We provide more specific results in the Appendix Figures 4, 5.

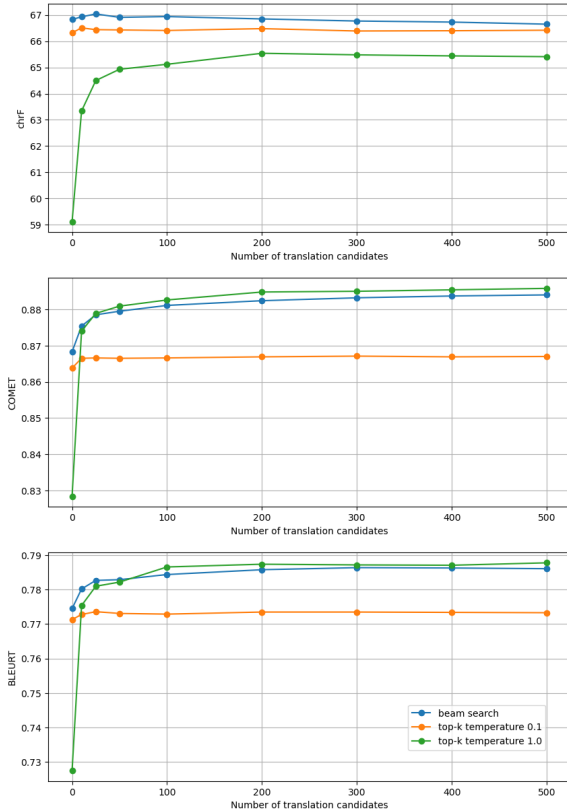


Figure 1: Comparison of beam search and top-k algorithms of the **Mix-tune** English–German model for the khresmoi test set. Top-k algorithm with temperature 1.0 showed superior performance on neural metrics over top-k with temperature 0.1 and slightly better performance than beam search. However, beam search achieved the highest score on the chrF metric, while the top-k algorithm with temperature 1.0 had the lowest score (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

Secondly, we decided to train our models on MBR-decoded data from 50 candidates selected by the beam search decoding algorithm. We considered the balance between improvement in evaluation metrics based on neural language models, stability across lexical metrics, and the execution time of MBR decoding, as shown in Figure 3.

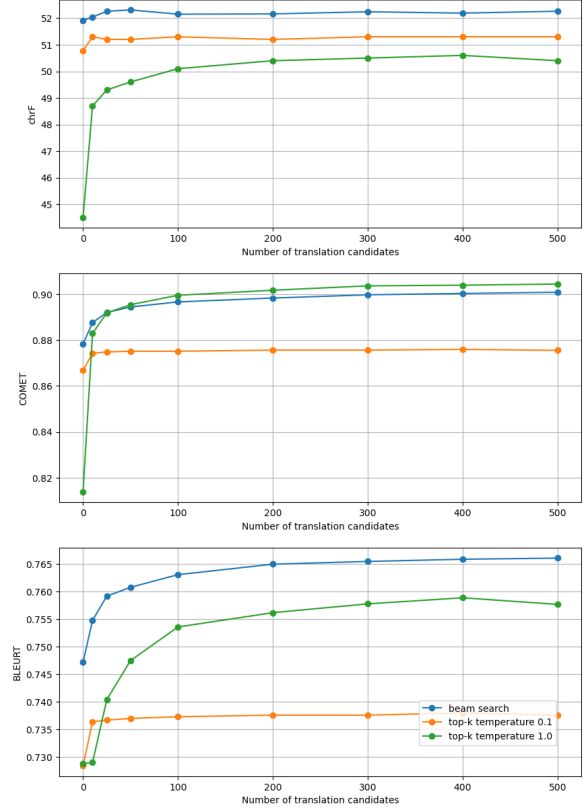


Figure 2: Comparison of beam search and top-k algorithms of the **baseline** Czech–Ukrainian model for the FLORES-200 test set. Beam search seems to be the superior option with the best performance on chrF and BLEURT metrics and slightly worse results on COMET over top-k with temperature 1.0 (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

We provide more detailed results in the Appendix Figures 6, 7, 8, 9, 10, 11, 12.

5.2 English to German

Table 1 shows the evaluation results on the in-domain test set khresmoi. All models self-improved with MBR decoding have shown enhanced performance. However, model **Mix-tune-domain-mbr-iter2** did not exhibit improvement over its first iteration **Mix-tune-domain-mbr**, even on COMET, which was the utility metric of MBR decoding. **Mix-tune-general-mbr** model shows a slightly better performance on BLEURT metric compared to models fine-tuned on in-domain MBR-decoded forward translations.

Table 2 presents the evaluation results on the FLORES-200 test set. Although chrF did not increase, the neural evaluation metrics showed improvement. Similar to the khresmoi test set, the **Mix-tune-domain-mbr-iter2** model showed a decrease in quality during the second iteration of self-

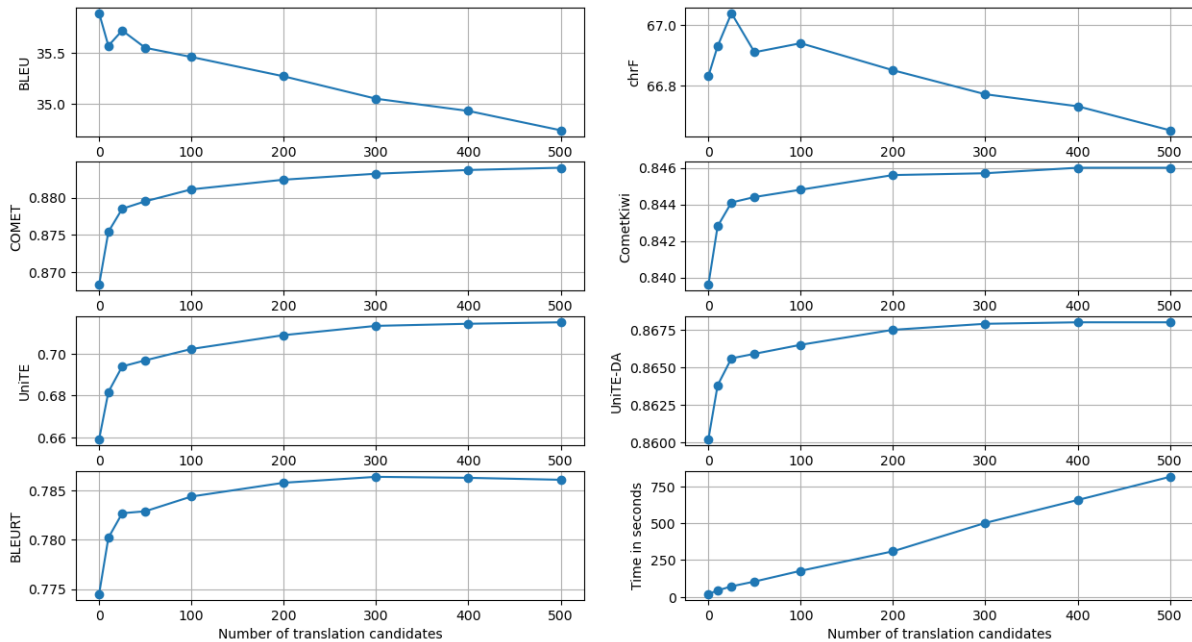


Figure 3: Comparison of beam search performance with a different number of samples of the **Mix-tune** English–German model for the khresmoi test set. Initial increases in the number of samples for MBR decoding showed very rapid gains, but further increases no longer resulted in such large gains, and performance on the n-gram metrics deteriorated (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

Model	chrF	COMET	BLEURT
Baseline	66.6	0.8653	0.7693
Mix-tune	66.8	0.8682	0.7749
Base-domain-mbr	66.9	0.8711*	0.7755
Mix-tune-domain-mbr	66.9	0.8728*	0.7792*
Mix-tune-domain-mbr-iter2	66.9	0.8727*	0.7791*
Mix-tune-general-mbr	66.9	0.8720*	0.7799*

Table 1: English–German khresmoi set results for the MBR self-improvement approaches. All models fine-tuned with MBR self-improvement technique have shown better performance over **Baseline** and **Mix-tune** models, including the **Mix-tune-general-mbr** model, which was finetuned on general-domain MBR-decoded data. The results marked with an asterisk (*) are statistically significant compared to the **Mix-tune** model.

improvement. **Mix-tune-general-mbr** showed superior performance over other models.

In summary, our findings demonstrate that applying MBR decoding significantly improves the performance of the high-resource English–German model for low-resource biomedical domain translation, particularly on neural network metrics. While lexical metrics show lower stability, they also hold potential for improvement.

Experiments demonstrated the robustness of self-improving models with the MBR decoding technique. Model fine-tuned on general forward translation had great performance on the in-domain test set and the model fine-tuned on

Model	chrF	COMET	BLEURT
Baseline	67.5	0.8751	0.7735
Mix-tune	67.5	0.8756	0.7744
Base-domain-mbr	67.2	0.8772	0.7743
Mix-tune-domain-mbr	67.3	0.8787*	0.7766
Mix-tune-domain-mbr-iter2	67.1	0.8766	0.7748
Mix-tune-general-mbr	67.5	0.8813*	0.7784*

Table 2: English–German FLORES-200 test set results for the MBR self-improvement approaches. **Mix-tune-general-mbr** model has shown superior performance, however, models with domain-specific forward translation maintain performance. The results marked with an asterisk (*) are statistically significant compared to the **Mix-tune** model.

domain-specific forward translation maintained performance on the general domain test set. We provide a broader evaluation in the Appendix Tables 9, 10, 11, 12.

5.3 Czech to Ukrainian

The results of the three MBR self-improvement approaches described in Section 4.6 are presented in Tables 3 and 4 for the FLORES-200 and WMT22 test sets, respectively.

We find that standard fine-tuning of the baseline model with MBR-decoded data yields the smallest improvements across all metrics, suggesting its limited effectiveness in this context. We note that both fine-tuning with a higher learning rate and

Model	chrF	COMET	BLEURT
Baseline	52.0	0.8779	0.7466
MBR-finetuned	52.4	0.8839	0.7522
MBR-ft-high-lr	52.7	0.8869	0.7553
MBR-resumed	52.7	0.8864	0.7557

Table 3: Czech–Ukrainian FLORES-200 test set results for the three MBR self-improvement approaches. All self-improved models exhibit improvements on all metrics compared to the baseline model, regardless of the fine-tuning approach used. Notably, both **MBR-ft-high-lr** and **MBR-resumed** models achieve the highest gains, demonstrating comparable performance. All self-improved models show statistical significance compared to the **Baseline** model.

Model	chrF	COMET	BLEURT
Baseline	58.4	0.8721	0.7498
MBR-finetuned	60.0	0.8803	0.7574
MBR-ft-high-lr	60.2	0.8844	0.7619
MBR-resumed	60.0	0.8852	0.7639

Table 4: Czech–Ukrainian WMT22 test set results for the three MBR self-improvement approaches. Similar to the FLORES-200 results, all self-improved models exhibit improvements on all metrics compared to the baseline model. However, on the WMT22 test set, the neural metrics favour the **MBR-resumed** model over the **MBR-ft-high-lr** model. All self-improved models show statistical significance compared to the **Baseline** model.

Model	chrF	COMET	BLEURT
Baseline	52.0	0.8779	0.7466
MBR-resumed	52.7*	0.8864*	0.7557*
MBR-resumed-iter2	52.8	0.8888*	0.7567
MBR-resumed-iter3	52.6	0.8901	0.7557

Table 5: Czech–Ukrainian iterative self-improvement results on the FLORES-200 test set. While the COMET score consistently improves across all three iterations, the chrF and BLEURT scores show a decrease in the third iteration. This suggests that the model overfits to COMET, harming the quality of the translation. Results with an asterisk (*) are statistically significant in comparison with the model in the row directly above it.

resuming the training exhibit comparable performance, with resumed training achieving slightly better results on the WMT22 test set. This may indicate that resuming training helps mitigate overfitting to the FLORES-200 validation set used during training.

Tables 5 and 6 showcase the impact of iterative training with MBR decoding on the FLORES-200 and WMT22 test sets, respectively. The second iteration consistently improves scores across all metrics, demonstrating the effectiveness of the

Model	chrF	COMET	BLEURT
Baseline	58.4	0.8721	0.7498
MBR-resumed	60.0*	0.8852*	0.7639*
MBR-resumed-iter2	60.3*	0.8885*	0.7641
MBR-resumed-iter3	60.1	0.8896	0.7578

Table 6: Czech–Ukrainian iterative self-improvement results on the WMT22 test set. Consistent with the FLORES-200 results, the COMET score improves across all iterations, while other metrics show a decrease in the last iteration. Notably, the BLEURT score not only decreases but falls below the score achieved by the first self-improved model. Results with an asterisk (*) are statistically significant in comparison with the model in the row directly above it.

iterative self-improvement process in refining the model’s translation capabilities. However, the third iteration leads to a decrease in both chrF and BLEURT scores. This suggests potential overfitting to the MBR decoding utility metric, where the model prioritizes aspects that score well according to COMET but may not translate to overall translation quality.

We provide extended evaluations in the Appendix in Tables 13, 14, 15, 16.

5.4 English to Hausa

Model	chrF	COMET	BLEURT	AfriCOMET
Baseline	49.9	0.7569	0.7931	0.6984
MBR-COMET	50.9	0.7720	0.8083	0.7207
MBR-AfriCOMET	51.2	0.7692	0.8061	0.7239

Table 7: English–Hausa FLORES-200 test set results for MBR self-improvement with different metrics. Both self-improved models achieve gains compared to the baseline model on all evaluation metrics. While the AfriCOMET-based model achieves a higher AfriCOMET score, reflecting its alignment with the specific evaluation metric, the COMET-based model surpasses it in both BLEURT and COMET scores, while showing a comparable gain on the AfriCOMET score. All self-improved models show statistical significance compared to the **Baseline** model.

Model	chrF	COMET	BLEURT	AfriCOMET
Baseline	51.6	0.7596	0.7791	0.6800
MBR-COMET	53.1	0.7752	0.7986	0.7046
MBR-AfriCOMET	53.0	0.7721	0.7956	0.7062

Table 8: English–Hausa NTREX test set results for MBR self-improvement with different metrics. Similar to the FLORES-200 results, both self-improved models using MBR decoding demonstrate improvements over the baseline model on all evaluation metrics. All self-improved models show statistical significance compared to the **Baseline** model.

This section compares the performance of

two MBR decoding self-improvement approaches for English–Hausa translation: one utilizing the WMT22 COMET model and another using the AfriCOMET model. The results are presented in Tables 7 and 8 for the FLORES-200 and NTREX test sets, respectively.

We observe that the AfriCOMET MBR-tuned model achieves gains over the WMT22 COMET MBR-tuned model on chrF for the FLORES-200 test set, but this advantage is not replicated on the NTREX test set. Additionally, the gains from AfriCOMET MBR-tuning are mainly limited to the AfriCOMET metric.

Our analysis reveals that the **MBR-AfriCOMET** model exhibits improvements over the **MBR-COMET** model primarily on lexical metrics in the case of the FLORES-200 test set, but not in the case of NTREX. The gains of the **MBR-AfriCOMET** model are mainly limited to AfriCOMET metrics, while other neural-based metrics consistently favour the **MBR-COMET** model.

While WMT22 COMET might exhibit a lower correlation with human judgment for the English–Hausa language pair than AfriCOMET, as reported by Wang et al. (2023), both self-improved models achieved significant and comparable gains on AfriCOMET. This suggests that WMT22 COMET, can still correctly rerank translation candidates and effectively guide the self-improvement process, leading to improvements on AfriCOMET, a metric specifically designed for African languages. This finding suggests that self-improvement guided by WMT22 COMET, with its diverse language coverage, might be effective even in low-resource settings, potentially reducing the need for additional adaptation of neural evaluation models to individual languages.

Additional evaluations are provided in the Appendix in Tables 17, 18.

6 Conclusion

This study demonstrated the effectiveness of model self-improvement through MBR decoding in improving translation quality. This approach proves beneficial for both high and low-resource languages, offering versatility in its application across diverse scenarios. Examples include domain-specific translation and the enhancement of general translation models.

We conducted experiments with various sample

sizes for MBR decoding, using two decoding algorithms: beam search and top-k. The aim was to find a balance between automatic metric gains and time efficiency. Our experiments have shown that the beam search algorithm with a beam size set to 50 is the optimal choice.

In the field of high-resource English-to-German biomedical translation, we investigated the impact of domain adaptation using various self-improvement approaches on MBR-decoded forward-translated data. Experiments showed that all MBR-based fine-tuning, regardless of the domain of the test set, improved performance compared to the baseline model. This finding highlights the robustness of the self-improvement technique.

Experiments on the Czech–Ukrainian language pair revealed that fine-tuning the MT model on MBR-decoded translations of the training data set significantly improves translation performance. Applying this process iteratively improves quality, but further iterations yield diminishing gains and at some point, the quality may even degrade due to overfitting to the MBR decoding utility metric.

In the English–Hausa experiments, we employed two models for MBR decoding: WMT22 COMET and AfriCOMET. Both models yielded comparable and significant improvements in automatic metrics, indicating their effectiveness in guiding the self-improvement process. While AfriCOMET, specifically trained on African languages, might intuitively seem favourable for this language pair, the performance of the **MBR-COMET** model highlights the potential of utilizing more widely applicable metrics like WMT22 COMET even for low-resource settings.

References

- Alabi, Jesujoba O., David Ifeoluwa Adelani, Marius Mosbach, and Dietrich Klakow. 2022. Adapting pre-trained language models to African languages via multilingual adaptive fine-tuning. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4336–4349, Gyeongju, Republic of Korea, October. International Committee on Computational Linguistics.
- Amrhein, Chantal and Rico Sennrich. 2022. Identifying weaknesses in machine translation metrics through minimum Bayes risk decoding: A case study for COMET. In He, Yulan, Heng Ji, Sujian Li, Yang Liu, and Chua-Hui Chang, editors, *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the*

- 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1125–1141, Online only, November. Association for Computational Linguistics.
- Bañón, Marta, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online, July. Association for Computational Linguistics.
- Burchell, Laurie, Alexandra Birch, Nikolay Bogoychev, and Kenneth Heafield. 2023. An open dataset and model for language identification. In Rogers, Anna, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 865–879, Toronto, Canada, July. Association for Computational Linguistics.
- Conneau, Alexis, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In Jurafsky, Dan, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online, July. Association for Computational Linguistics.
- Dušek, Ondřej, Jan Hajič, Jaroslava Hlaváčová, Jindřich Libovický, Pavel Pecina, Aleš Tamchyna, and Zdeňka Urešová. 2017. Khresmoi summary translation test data 2.0. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.
- Fernandes, Patrick, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In Carpuat, Marine, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States, July. Association for Computational Linguistics.
- Finkelstein, Mara, Subhajt Naskar, Mehdi Mirzazadeh, Apurva Shah, and Markus Freitag. 2023. Mbr and qe finetuning: Training-time distillation of the best and most expensive decoding methods.
- Freitag, Markus, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.
- Freitag, Markus, Behrooz Ghorbani, and Patrick Fernandes. 2023. Epsilon sampling rocks: Investigating sampling strategies for minimum bayes risk decoding for machine translation. *arXiv preprint arXiv:2305.09860*.
- Gowda, Thamme, Zhao Zhang, Chris Mattmann, and Jonathan May. 2021. Many-to-English machine translation tools, data, and pretrained models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 306–316, Online, August. Association for Computational Linguistics.
- Gulcehre, Caglar, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, et al. 2023. Reinforced self-training (rest) for language modeling. *arXiv preprint arXiv:2308.08998*.
- Jon, Josef, Martin Popel, and Ondřej Bojar. 2023. CUNI at WMT23 general translation task: MT and a genetic algorithm. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 119–127, Singapore, December. Association for Computational Linguistics.
- Junczys-Dowmunt, Marcin, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia, July. Association for Computational Linguistics.
- Kocmi, Tom, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore, December. Association for Computational Linguistics.
- Koehn, Philipp. 2004. Statistical significance tests for machine translation evaluation. In Lin, Dekang and Dekai Wu, editors, *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain, July. Association for Computational Linguistics.

- Kudo, Taku and John Richardson. 2018. Sentence-Piece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In Blanco, Eduardo and Wei Lu, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium, November. Association for Computational Linguistics.
- Kudo, Keito, Takumi Ito, Makoto Morishita, and Jun Suzuki. 2023. SKIM at WMT 2023 general translation task. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 128–136, Singapore, December. Association for Computational Linguistics.
- Kumar, Shankar and William Byrne. 2004. Minimum Bayes-risk decoding for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 169–176, Boston, Massachusetts, USA, May 2 - May 7. Association for Computational Linguistics.
- Neves, Mariana, Antonio Jimeno Yepes, and Aurélie Névéol. 2016. The scielo corpus: a parallel corpus of scientific publications for biomedicine. In Calzolari, Nicoletta, Khalid Choukri, Thierry Declerck, Sara Goggi, Marko Grobelnik, Bente Maegaard, Joseph Mariani, Helene Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 2942–2948, Portorož, Slovenia, May. European Language Resources Association (ELRA).
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Searley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.
- Nowakowski, Artur, Gabriela Pałka, Kamil Guttman, and Mikołaj Pokrywka. 2022. Adam Mickiewicz University at WMT 2022: NER-assisted and quality-aware neural machine translation. In Koehn, Philipp, Loic Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Tom Kocmi, Andre Martins, Makoto Morishita, and Christof Monz, editors, *Proceedings of the Sixth Conference on Machine Translation*, pages 1030–1040, Online, November. Association for Computational Linguistics.
- Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 326–334, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Och, Franz Josef. 2003. Minimum error rate training in statistical machine translation. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics*, pages 160–167, Sapporo, Japan, July. Association for Computational Linguistics.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Isabelle, Pierre, Eugene Charniak, and Dekang Lin, editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Popović, Maja. 2015. chrF: character n-gram F-score for automatic MT evaluation. In Bojar, Ondřej, Rajan Chatterjee, Christian Federmann, Barry Haddow, Chris Hokamp, Matthias Huck, Varvara Logacheva, and Pavel Pecina, editors, *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal, September. Association for Computational Linguistics.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels, October. Association for Computational Linguistics.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In Webber, Bonnie, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.
- Rei, Ricardo, Ana C Farinha, Chrysoula Zerva, Daan van Stigt, Craig Stewart, Pedro Ramos, Taisiya Glushkova, André F. T. Martins, and Alon Lavie. 2021. Are references really needed? unbabel-IST 2021 submission for the metrics shared task. In Barrault, Loic, Ondrej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Tom Kocmi, Andre Martins, Makoto Morishita, and Christof Monz, editors, *Proceedings of the Sixth Conference on Machine Translation*, pages 1030–1040, Online, November. Association for Computational Linguistics.

- Rei, Ricardo, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022. CometKiwI: IST-unbabel 2022 submission for the quality estimation shared task. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Rei, Ricardo, Nuno M. Guerreiro, Marcos Treviso, Luisa Coheur, Alon Lavie, and André Martins. 2023. The inside story: Towards better understanding of machine translation neural evaluation metrics. In Rogers, Anna, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1089–1105, Toronto, Canada, July. Association for Computational Linguistics.
- Sellam, Thibault, Dipanjan Das, and Ankur Parikh. 2020a. BLEURT: Learning robust metrics for text generation. In Jurafsky, Dan, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online, July. Association for Computational Linguistics.
- Sellam, Thibault, Dipanjan Das, and Ankur P Parikh. 2020b. Bleurt: Learning robust metrics for text generation. In *Proceedings of ACL*.
- Tiedemann, Jörg and Lars Nygaard. 2004. The OPUS corpus - parallel and free: <http://logos.uio.no/opus>. In Lino, Maria Teresa, Maria Francisca Xavier, Fátima Ferreira, Rute Costa, and Raquel Silva, editors, *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)*, Lisbon, Portugal, May. European Language Resources Association (ELRA).
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Guyon, I., U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Vernikos, Giorgos and Andrei Popescu-Belis. 2024. Don't rank, combine! combining machine translation hypotheses using quality estimation.
- Wan, Yu, Dayiheng Liu, Baosong Yang, Haibo Zhang, Boxing Chen, Derek Wong, and Lidia Chao. 2022. UniTE: Unified translation evaluation. In Muresan, Smaranda, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8117–8127, Dublin, Ireland, May. Association for Computational Linguistics.
- Wang, Jiayi, David Ifeoluwa Adelani, Sweta Agrawal, Ricardo Rei, Eleftheria Briakou, Marine Carpuat, Marek Masiak, Xuanli He, Sofia Bourhim, Andiswa Bukula, Muhidin Mohamed, Temitayo Olatoye, Hamam Mokayed, Christine Mwase, Wangui Kimotho, Foutse Yuehgo, Anuoluwapo Aremu, Jessica Ojo, Shamsuddeen Hassan Muhammad, Salomey Osei, Abdul-Hakeem Omotayo, Chiamaka Chukwunke, Perez Ogayo, Oumaima Hourrane, Salma El Anigri, Lolwethu Ndolela, Thabiso Mangwana, Shafie Abdi Mohamed, Ayinde Hassan, Oluwabusayo Olufunke Awoyomi, Lama Alkhaled, Sana Al-Azzawi, Naome A. Etori, Millicent Ochieng, Clemencia Siro, Samuel Njoroge, Eric Muchiri, Wangari Kimotho, Lyse Naomi Wamba Momo, Daud Abolade, Simbiat Ajao, Tosin Adewumi, Iyanuoluwa Shode, Ricky Macharm, Ruqayya Nasir Iro, Saheed S. Abdullahi, Stephen E. Moore, Bernard Opoku, Zainab Akinjobi, Abeeb Afolabi, Nnaemeka Obiefuna, Onyekachi Raphael Ogbu, Sam Brian, Verrah Akinyi Otiende, Chinedu Emmanuel Mbonu, Sakayo Toadoum Sari, and Pontus Stenetorp. 2023. Afrimte and africomet: Empowering comet to embrace under-resourced african languages.
- Zaragoza-Bernabeu, Jaume, Gema Ramírez-Sánchez, Marta Bañón, and Sergio Ortiz Rojas. 2022. Bicleaner AI: Bicleaner goes neural. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 824–831, Marseille, France, June. European Language Resources Association.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT
Baseline	66.6	35.6	0.8653	0.8373	0.6441	0.8574	0.7693
Mix-tune	66.8	35.9	0.8682	0.8397	0.6594	0.8602	0.7749
Base-domain-mbr	66.9	35.7	0.8711	0.8416	0.6694	0.8621	0.7755
Mix-tune-domain-mbr	66.9	35.8	0.8728	0.8423	0.6766	0.8631	0.7792
Mix-tune-domain-mbr-iter2	66.9	35.6	0.8727	0.8423	0.6757	0.8633	0.7791
Mix-tune-general-mbr	66.9	35.5	0.8720	0.8422	0.6775	0.8631	0.7799

Table 9: English–German khresmoi set results for the MBR self-improvement approaches. All models fine-tuned with MBR self-improvement technique have shown better performance over **Baseline** and **Mix-tune** models, even **Mix-tune-general-mbr** model with general forward translations.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT
Baseline	63.1	35.0	0.8505	0.8336	0.5368	0.8470	0.7500
Mix-tune	63.5	35.6	0.8525	0.8360	0.5418	0.8495	0.7541
Base-domain-mbr	63.5	35.8	0.8549	0.8374	0.5549	0.8501	0.7522
Mix-tune-domain-mbr	63.6	35.7	0.8540	0.8379	0.5552	0.8508	0.7530
Mix-tune-domain-mbr-iter2	63.7	35.9	0.8543	0.8383	0.5575	0.8510	0.7535
Mix-tune-general-mbr	63.4	35.4	0.8547	0.8378	0.5613	0.8501	0.7542

Table 10: English–German WMT22-medline set results for the MBR self-improvement approaches. All models fine-tuned with MBR self-improvement technique have shown better performance over **Mix-tune** model except on metric BLEURT. On this specific test set, **Mix-tune-domain-mbr-iter2** outperformed the **Mix-tune-domain-mbr** model, unlike the results observed on other test sets.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT
Baseline	67.5	42.0	0.8751	0.8454	0.6630	0.8614	0.7735
Mix-tune	67.5	42.2	0.8756	0.8457	0.6657	0.8617	0.7744
Base-domain-mbr	67.2	41.7	0.8772	0.8469	0.6677	0.8632	0.7743
Mix-tune-domain-mbr	67.3	41.7	0.8787	0.8477	0.6719	0.8641	0.7766
Mix-tune-domain-mbr-iter2	67.1	41.5	0.8766	0.8466	0.6653	0.8629	0.7748
Mix-tune-general-mbr	67.5	41.8	0.8813	0.8484	0.6824	0.8654	0.7784

Table 11: English–German FLORES-200 test set results for the MBR self-improvement approaches. **Mix-tune-general-mbr** model has shown superior performance, however, models with domain-specific forward translation maintain performance.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT
Baseline	63.8	36.6	0.8428	0.8328	0.5308	0.8420	0.7106
Mix-tune	63.7	36.5	0.8427	0.8322	0.5283	0.8414	0.7107
Base-domain-mbr	63.3	35.8	0.8463	0.8359	0.5376	0.8454	0.7138
Mix-tune-domain-mbr	63.2	35.9	0.8468	0.8358	0.5404	0.8464	0.7132
Mix-tune-domain-mbr-iter2	63.0	35.5	0.8460	0.8345	0.5348	0.8455	0.7119
Mix-tune-general-mbr	64.1	36.7	0.8629	0.8399	0.5622	0.8492	0.7202

Table 12: English–German Statmt test set results for the MBR self-improvement approaches. **Mix-tune-general-mbr** model has shown significantly improved performance on every metric, however models with domain-specific forward translation maintain performance.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT
Baseline	52.0	22.2	0.8779	0.8449	0.4441	0.9017	0.7466
MBR-finetuned	52.4	22.3	0.8839	0.8513	0.4715	0.9063	0.7522
MBR-ft-high-lr	52.7	22.6	0.8869	0.8543	0.4829	0.9085	0.7553
MBR-resumed	52.7	22.8	0.8864	0.8540	0.4824	0.9086	0.7557

Table 13: Extended Czech–Ukrainian FLORES-200 test set results for the three MBR self-improvement approaches. All approaches lead to an increase in evaluation scores. Both **MBR-ft-high-lr** and **MBR-resumed** models achieve the highest gains, demonstrating comparable performance.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT
Baseline	58.4	31.1	0.8721	0.8046	0.3744	0.8795	0.7498
MBR-finetuned	60.0	32.3	0.8803	0.8121	0.4112	0.8846	0.7574
MBR-ft-high-lr	60.2	33.2	0.8844	0.8152	0.4246	0.8880	0.7619
MBR-resumed	60.0	33.0	0.8852	0.8162	0.4236	0.8890	0.7639

Table 14: Extended Czech–Ukrainian WMT22 test set results for the three MBR self-improvement approaches. As in the case of evaluation results on the FLORES-200 test set, all approaches improve upon the baseline model, although **MBR-resumed** stands out across all neural metrics apart from UniTE.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT
Baseline	52.0	22.2	0.8779	0.8449	0.4441	0.9017	0.7466
MBR-resumed	52.7	22.8	0.8864	0.8540	0.4824	0.9086	0.7557
MBR-resumed-iter2	52.8	22.6	0.8888	0.8557	0.4882	0.9099	0.7567
MBR-resumed-iter3	52.6	22.3	0.8901	0.8562	0.4873	0.9097	0.7557

Table 15: Extended Czech–Ukrainian iterative self-improvement results on the FLORES-200 test set. Models increase in quality across all neural metrics until the third iteration, when the quality measured by metrics other than COMET and CometKiwi decreases. It’s worth noticing that the BLEU score increases only in the first iteration and slowly degrades in consecutive iterations.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT
Baseline	58.4	31.1	0.8721	0.8046	0.3744	0.8795	0.7498
MBR-resumed	60.0	33.0	0.8852	0.8162	0.4236	0.8890	0.7639
MBR-resumed-iter2	60.3	32.6	0.8885	0.8183	0.4349	0.8900	0.7641
MBR-resumed-iter3	60.1	31.9	0.8896	0.8174	0.4312	0.8887	0.7578

Table 16: Extended Czech–Ukrainian iterative self-improvement results on the WMT22 test set. Evaluations across all metrics show similar tendencies as in the case of FLORES-200, except for CometKiwi which also decreases in the third iteration.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT	AfriCOMET
Baseline	49.9	22.3	0.7569	0.5597	-0.2297	0.6082	0.7931	0.6984
MBR-COMET	50.9	23.2	0.7720	0.5707	-0.1777	0.6233	0.8083	0.7207
MBR-AfriCOMET	51.2	23.4	0.7692	0.5638	-0.1878	0.6183	0.8061	0.7239

Table 17: Extended English–Hausa results on the FLORES-200 test set. According to lexical metrics and AfriCOMET, the **MBR-AfriCOMET** model shows the greatest improvement. However, other neural metrics suggest that the **MBR-COMET** model is superior.

Model	chrF	BLEU	COMET	CometKiwi	UniTE	UniTE-DA	BLEURT	AfriCOMET
Baseline	51.6	23.9	0.7596	0.5704	-0.1763	0.6294	0.7791	0.6800
MBR-COMET	53.1	25.3	0.7752	0.5865	-0.1051	0.6484	0.7986	0.7046
MBR-AfriCOMET	53.0	24.9	0.7721	0.5803	-0.1273	0.6409	0.7956	0.7062

Table 18: Extended English–Hausa results on the NTREX test set. In contrast to evaluations on the FLORES-200 test set, in this case only the AfriCOMET metric favours the **MBR-AfriCOMET** model.

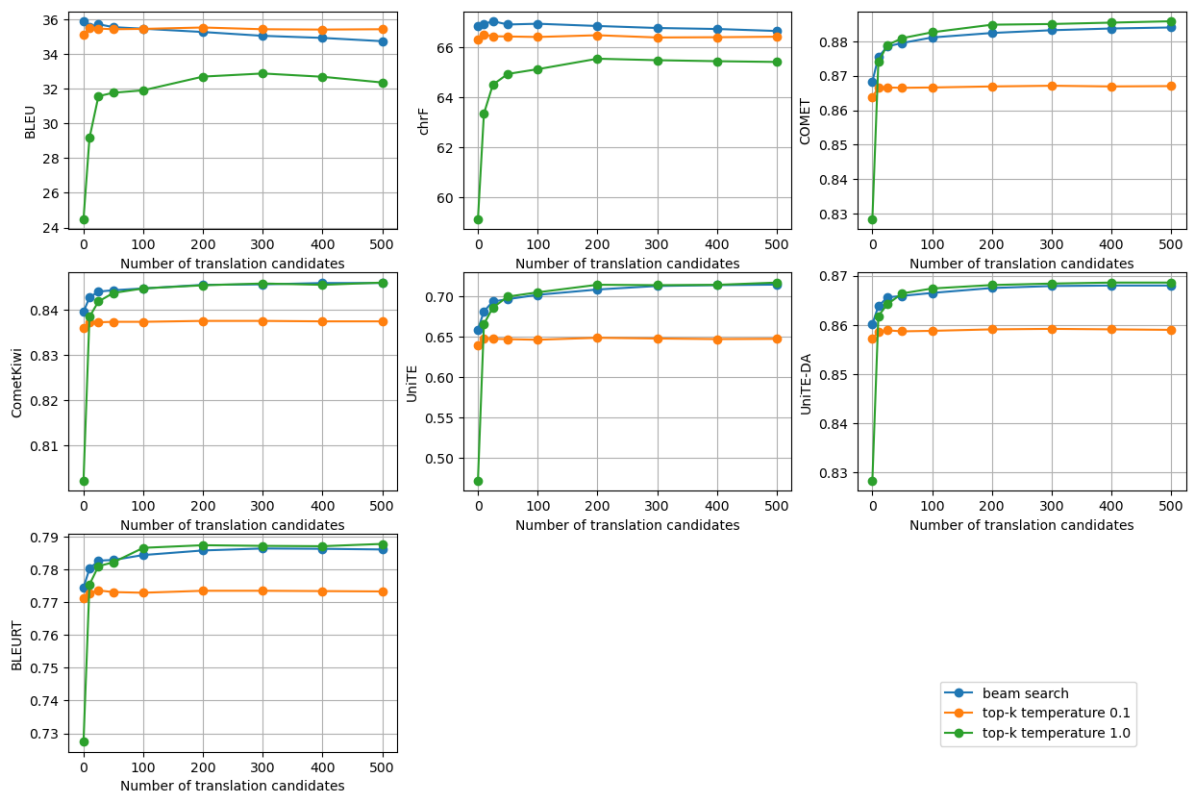


Figure 4: Comparison of beam search and top-k algorithms of the **Mix-tune** English–German model for the khresmoi test set. Top-k algorithm with temperature 1.0 showed superior performance on neural metrics over top-k with temperature 0.1 and slightly better performance than beam search. However, beam search achieved the highest score on the chrF metric, while the top-k algorithm with temperature 1.0 had the lowest score for lexical metrics (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

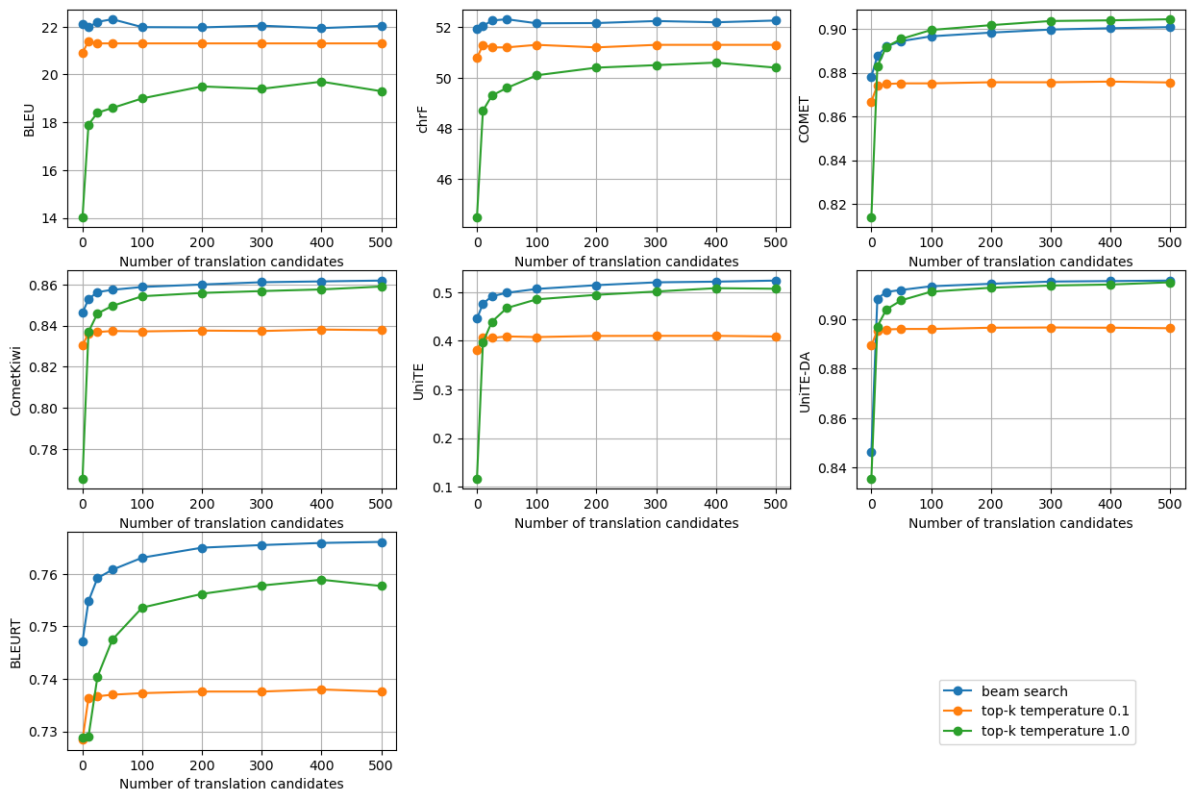


Figure 5: Comparison of beam search and top-k algorithms of the **baseline** Czech–Ukrainian model for the FLORES-200 test set. Beam search seems to be the superior option with the best performance on every metric except COMET (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

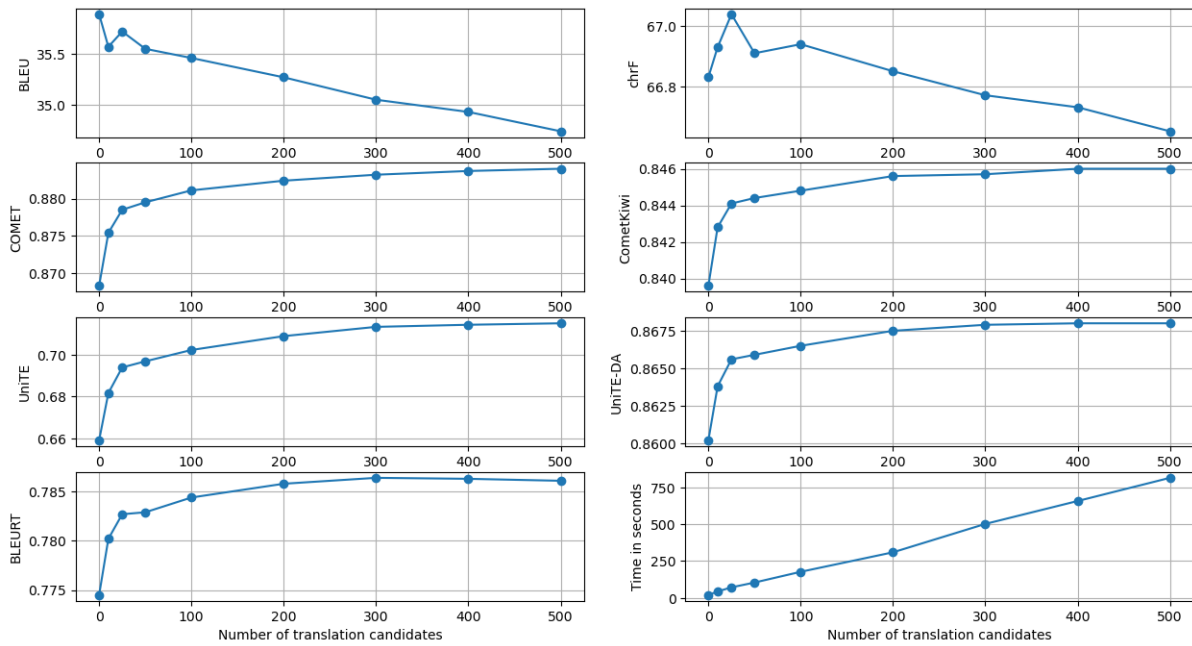


Figure 6: Comparison of beam search performance with a different number of samples of the **Mix-tune** English–German model for the khresmoi test set. Initial increases in the number of samples for MBR decoding showed very rapid gains, but further increases no longer resulted in such large gains and performance on the n-gram metrics deteriorated (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

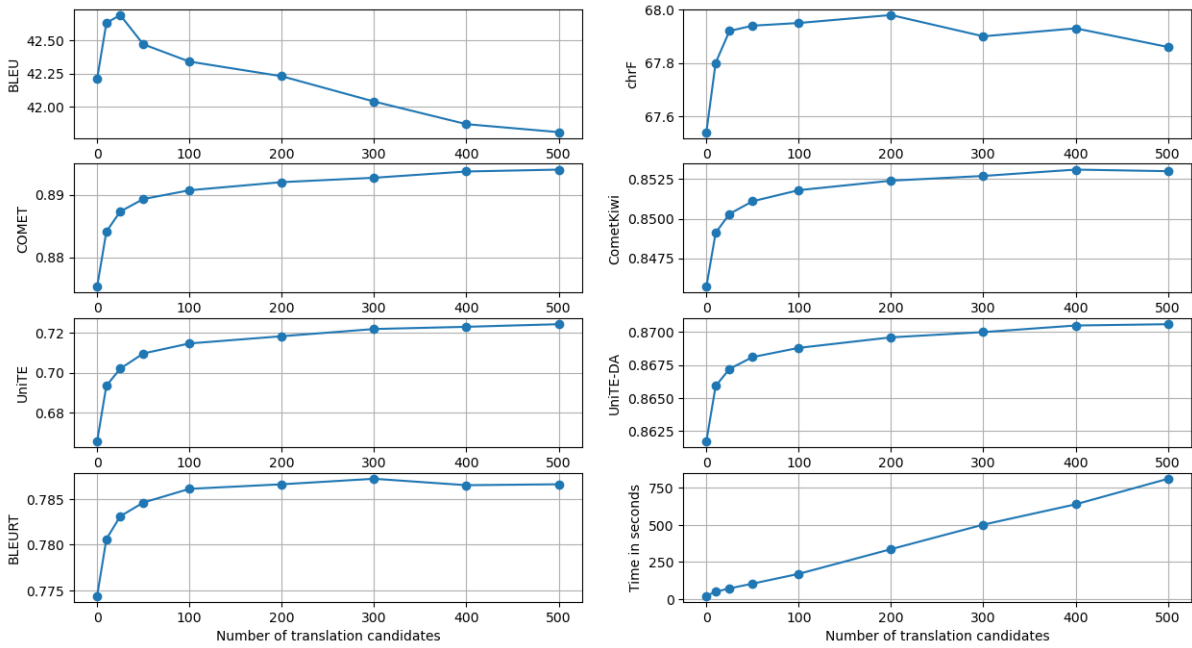


Figure 7: Comparison of beam search performance with a different number of samples of the **Mix-tune** English–German model for the FLORES-200 test set. Initial increases in the number of samples for MBR decoding showed very rapid gains, but further increases no longer resulted in such large gains and performance on the n-gram metrics deteriorated (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

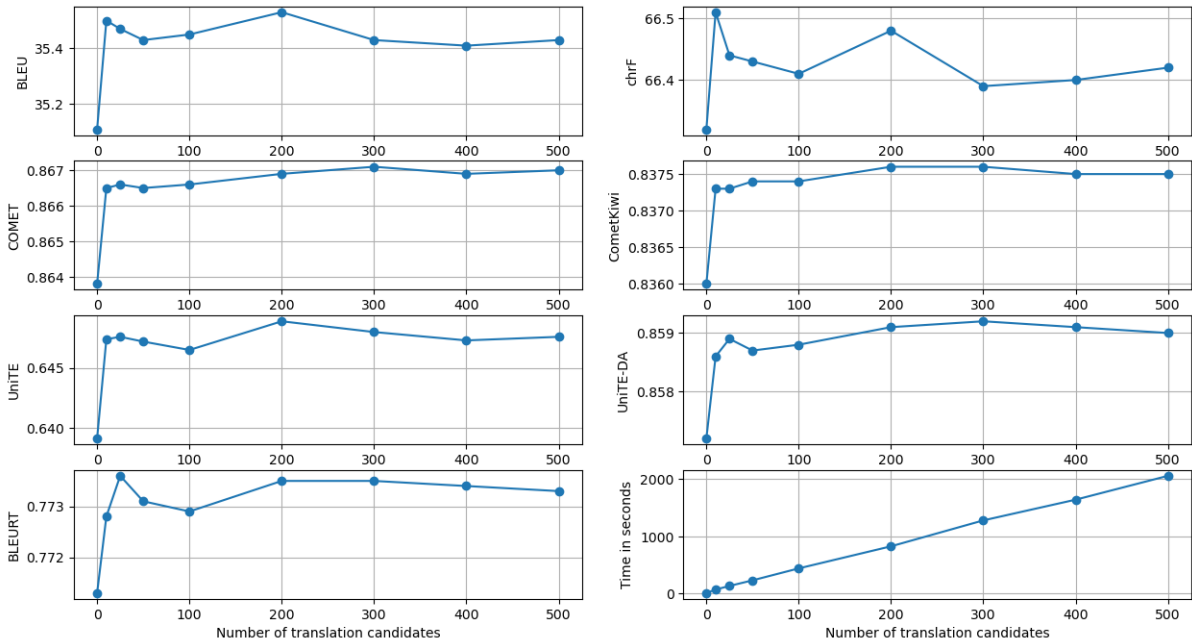


Figure 8: Comparison of top-k performance (temperature 0.1, k=10) with different number of samples of the **Mix-tune** English–German model for the khresmoi test set. Initial increases in the number of samples for MBR decoding showed very rapid gains, but further increases no longer resulted in such large gains (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

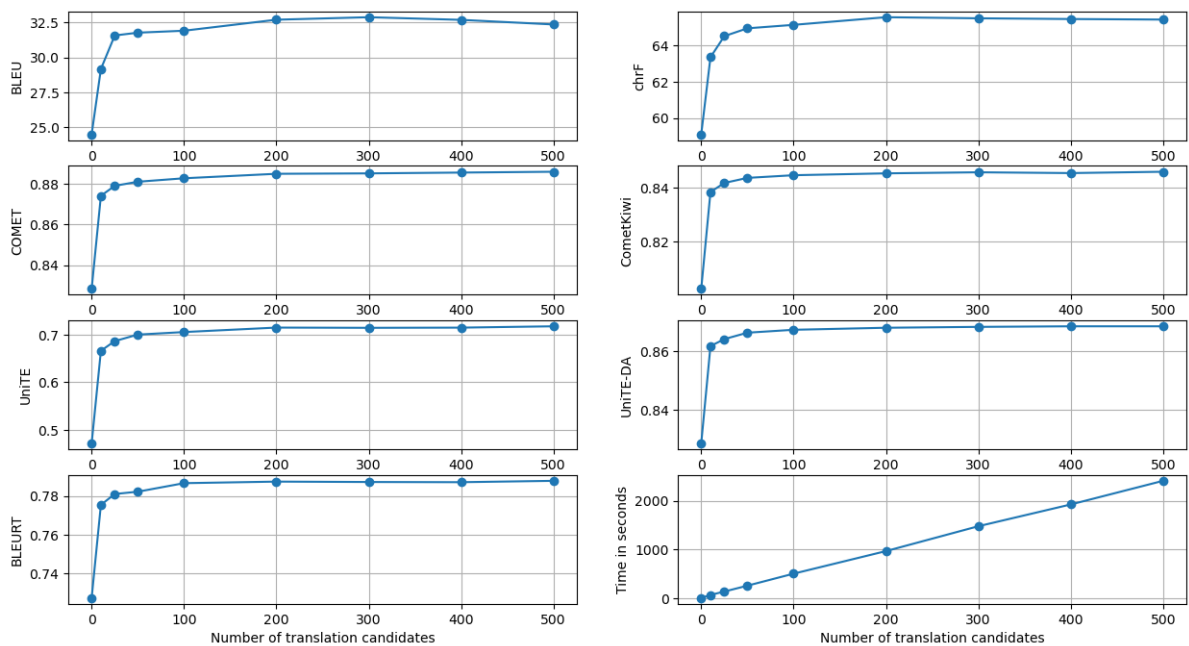


Figure 9: Comparison of top-k performance (temperature 1.0, k=10) with different number of samples of the **Mix-tune** English–German model for the khresmoi test set. Initial increases in the number of samples for MBR decoding showed very rapid gains, but further increases no longer resulted in such large gains (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

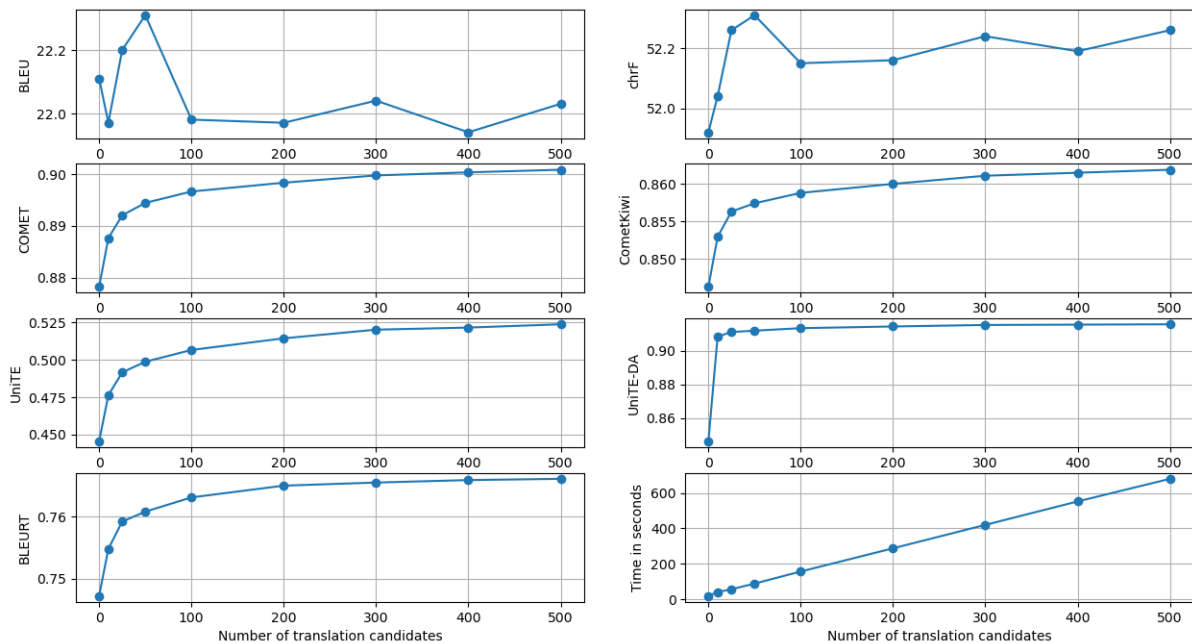


Figure 10: Comparison of beam search performance with different number of samples of the **Baseline** Czech–Ukrainian model for the FLORES-200 test set. Initial increases in the number of samples for MBR decoding showed very rapid gains, but further increases no longer resulted in such large gains and performance on the n-gram metrics deteriorated (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

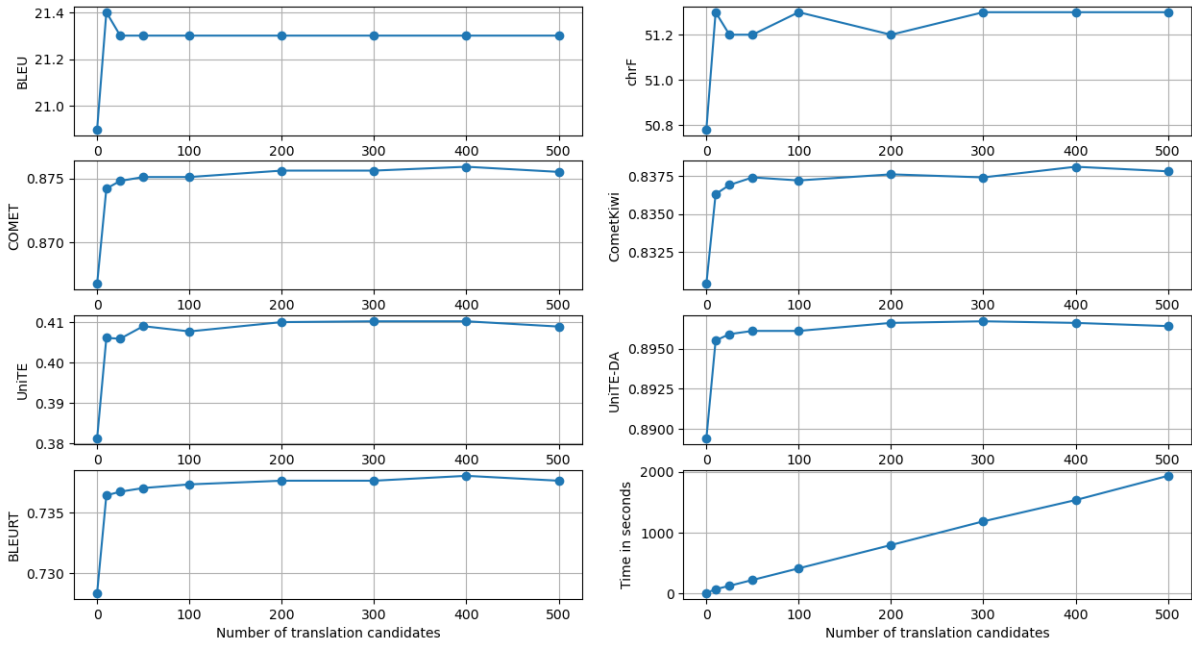


Figure 11: Comparison of top-k performance (temperature 0.1, $k=10$) with different number of samples of the **Baseline** Czech-Ukrainian model for the FLORES-200 test set. Initial increases in the number of samples for MBR decoding showed very rapid gains, but further increases no longer resulted in such large gains (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

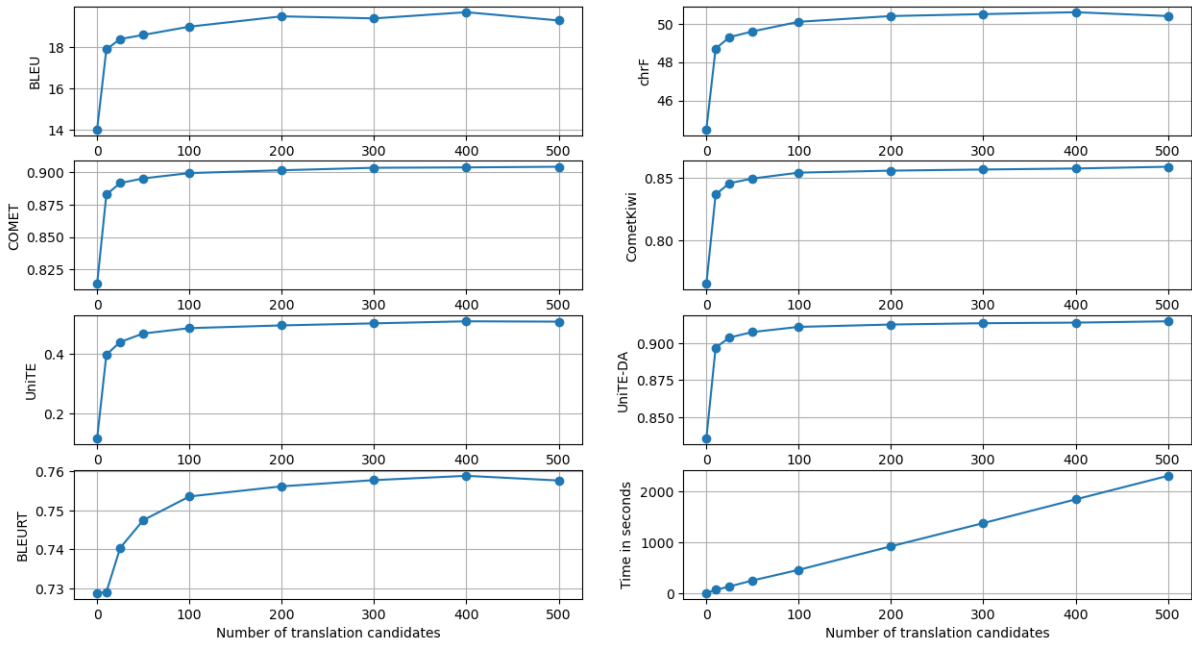


Figure 12: Comparison of top-k performance (temperature 1.0, $k=10$) with different number of samples of the **Baseline** Czech-Ukrainian model for the FLORES-200 test set. Initial increases in the number of samples for MBR decoding showed very rapid gains, but further increases no longer resulted in such large gains (translation without MBR decoding is represented on the chart as the number of translation candidates equal to 0).

Mitra: Improving Terminologically Constrained Translation Quality with Backtranslations and Flag Diacritics

Iikka Hauhio^{†‡} and Théo Friberg^{†‡}

[†] Kielikone Oy, Helsinki, Finland

[‡] Department of Computer Science, University of Helsinki, Finland

{iikka.hauhio,theo.friberg}@kielikone.fi

Abstract

Terminologically constrained machine translation is a hot topic in the field of neural machine translation. One major way to categorize constrained translation methods is to divide them into “hard” constraints that are forced into the target language sentence using a special decoding algorithm, and “soft” constraints that are included in the input given to the model.

We present a constrained translation pipeline that combines soft and hard constraints while being completely model-agnostic, i.e. our method can be used with any NMT or LLM model. In the “soft” part, we substitute the source language terms in the input sentence for the backtranslations of their target language equivalents. This causes the source sentence to be more similar to the intended translation, thus making it easier to translate for the model. In the “hard” part, we use a novel nondeterministic finite state transducer-based (NDFST) constraint recognition algorithm utilizing flag diacritics to force the model to use the desired target language terms.

We test our model with both Finnish–English and English–Finnish real-world vocabularies. We find that our methods consistently improve the translation quality when compared to previous constrained decoding algorithms, while the improvement over unconstrained translations depends on the

familiarity of the model over the subject vocabulary and the quality of the vocabulary.

1 Introduction

In this paper, we present *Mitra*, an end-to-end pipeline for terminology-constrained translation that combines a novel constrained beam search algorithm with backtranslation substitution.

Terminology-constrained machine translation is a popular topic in the field of machine translation, and has been a focus of several shared tasks in the WMT conference (Alam et al. 2021b; Semenov et al. 2023). In constrained translation, the system is given a lexicon, or a terminology, and it must use the words given in this terminology when translating sentences. While this was a trivial task in phrase-based statistical (cf. Koehn et al. 2003) and rule-based machine translation systems (cf. Arnola 1996), implementing it for neural systems has proved to be much more difficult due to their black-box nature.

The existing methods can be divided into the so called “hard” and “soft” constraints. Hard constraints use *constrained decoding* algorithms such as constrained beam-search (Hokamp and Liu 2017; Anderson et al. 2017), which first decides on the acceptable forms of constraints at the token level and then forces the decoder of an NMT system to abide by them. Soft methods, on the other hand, use a neural network specifically trained for the purpose of constrained translation, and the constraints can be given to the encoder of the network as input (cf. Bergmanis and Pinnis 2021). Both of these methods have their own advantages: hard constraints can be enforced on any neural network without the need to train or fine-tune anything, while soft constraints are generally faster. Typically, hard constraints only mandate that terms occur somewhere in the trans-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

lated sentence, while soft constraint approaches allow an explicit coupling of the source and target terms.

Implementing hard constraints for agglutinative languages has several difficulties. Foremost, if a term has multiple possible inflected forms, and the correct form is not known beforehand, the constrained decoding algorithm must be given multiple alternative forms (Anderson et al. 2017), a feature not widely supported by many algorithms such as Post and Vilar (2018); Hu et al. (2019). Moreover, the system must be able to generate these alternative forms, requiring the use of language-specific morphological generators, which might not be available for all languages.

Another problem of hard constraints is that the translation quality might be very poor if the machine translation model does not recognize the constrained words, a situation which in our experience is very common as terminologies often contain uncommon technical jargon, brand names, and other terms not appearing in the training data. In our specific case, we tested our method with a vocabulary provided by the Finnish Forest Centre containing names of insects in Finnish and English (Metsäkeskus 2023). Many of the names have surprising translations: for example, a “violet tan-bark beetle” is called “papintappaja” in Finnish, which means “priest-slayer” if translated literally. If the translation model has not seen this term before, it cannot correctly translate it without terminology constraints. Even then, the model has a hard time determining the correct location for the constraint in the output sentence (cf. Hasler et al. 2018). See Section 2.3 for more details of this problem.

We propose a combined method that tackles both of the aforementioned problems. To support heavily-inflected languages such as Finnish, we introduce a “hard” finite-state automaton-based constraint recognition algorithm that can recognize arbitrarily large disjunctive constraints. For the problem of expressions that were not encountered at training time, we propose a “soft” backtranslation substitution algorithm that makes it possible to use terminology constraints even when the neural network sees no connection between the source-language term and the target-language term. Both of these methods are integrated into an end-to-end pipeline that takes source-language sentences and lexicons as input and produces lexically constrained translations. We argue that the “hard” and

“soft” constraint methods complement each other and work together as a whole greater than the sum of its parts.

In this paper, we first describe existing constrained beam search algorithms (Section 2). We then give an overview of our pipeline, including detailed descriptions of the backtranslation and the constraint recognition algorithms (Section 3). Finally, we evaluate these algorithms against the existing algorithms (Section 4).

2 Constrained Beam Search

“Hard” terminology constraints refer to phrases (i.e. sequences of tokens) that are forced to appear in an output sequence during beam search. While a regular beam search compiles the list of new hypotheses by finding the most probable continuations for the current hypotheses (Koehn 2009), a constrained beam search algorithm additionally proposes tokens in the constraints as possible continuations (Hokamp and Liu 2017). Several algorithms exist, differing mainly in beam allocation (i.e. how much of the beam is reserved for hypotheses containing constraints), and constraint recognition (i.e. how they determine which constraints are fulfilled and propose constraint tokens as continuations).

2.1 Existing Algorithms

Hokamp and Liu (2017) present an algorithm called Grid Beam Search (GBS), which allocates $C + 1$ hypothesis banks in the beam, where C is the number of constraint tokens. Each hypothesis in bank $i \in [0, C]$ must have exactly i fulfilled constraint tokens. Unlike all other algorithms inspected, they do not allow backtracking, i.e. constraints that were previously considered fulfilled can not become unfulfilled again: if a model begins generating a prefix of a multi-token constraint, the hypothesis can only be continued by generating the rest of the constraint. In other words, hypotheses can only move upwards in the banks or stay at the same level. This also means that GBS does not require a constraint recognition algorithm for detecting which constraints are currently fulfilled in a hypothesis – it only needs to remember which constraint tokens it has previously generated, since those tokens will stay fulfilled.

Post and Vilar (2018) criticize GBS for its beam allocation, as the number of hypotheses grows linearly with the number of constraints. They propose a method called Dynamic Beam Allocation (DBA), in which the beam size is constant, and the hypothe-

Original sentence	Metsäpaloregimi summaa yhteen lähes kaikki metsäpaloihin vaikuttavat tekijät.
Greedy tracking	The <u>forest fire</u> [†] forest fire regime brings together almost all the factors that affect forest fires.
Exact tracking	The forest fire regime brings together almost all the factors affecting <u>forest fires</u> .

Table 1: Having reached the point marked with [†], greedy tracking accepts the constraint `metsäpalo → forest fire` and discards the start of the constraint `metsäpaloregimi → forest fire regime`. See also Appendix D.1.1.

Original sentence	[...] rajoittamaan ja supistamaan paloa rajoituslinjojen avulla.
Natural translation	[...] limit and reduce the fire by means of firebreaks .
Machine translation	[...] limit and reduce the firebreak by means <u>limiting lines</u>
Original sentence	Raivaamalla tehtyjä rajoituslinjoja ₁ ovat palokuja ₂ ja palokäytävä ₃ .
Natural translation	Fire lines ₂ and fire alleys ₃ are firebreaks ₁ made by clearing.
Machine translation	The firebreaks ₁ and fire alley ₃ are fire lines ₂ made by clearing.

Table 2: Examples of leaks and misplacements. In the first sentence, the model leaks the structure of “rajoituslinja” (lit. limiting line). In the second sentence, the model exchanges the constraints, changing the meaning of the sentence.

ses are assigned to different numbers of fulfilled constraints dynamically, making GPU memory optimization easier. In addition, to enable backtracking, they give a detailed description of a table-based data structure used for constraint recognition. The table contains information on which tokens are part of multi-token constraints and which of them are fulfilled. If the algorithm generates a token that is not a continuation of the current constraint being generated, it backtracks by marking the previously generated tokens unfulfilled.

This algorithm is further improved by Hu et al. (2019) who note that the constraint recognition algorithm proposed by Post and Vilar is flawed, as it cannot properly recognize overlapping constraints. They propose a trie-based algorithm claimed to resolve these issues. They also detail a method that allows sorting and selecting hypotheses completely in the GPU memory, further decreasing the overhead of the algorithm. Neither Post and Vilar (2018) nor Hu et al. (2019) support disjunctive constraints required by heavily-inflected languages, although we note that either of the algorithms can be relatively easily expanded to support them.

In addition to tables and tries, finite-state automata can be used to recognize fulfilled constraints in a hypothesis (Anderson et al. 2017; Hasler et al. 2018). This approach also supports disjunctive constraints and multi-token constraints.

2.2 The Problem of Greediness

The constraint recognition algorithms proposed by Post and Vilar (2018) and Hu et al. (2019) are greedy, which allows them to operate in $O(n)$ time. While this is good for time complexity, it also

makes the algorithms incorrect: they cannot detect some valid sequences that contain overlapping constraints (see Table 1 for an example).

We assert that no greedy algorithm can detect all valid sequences. Consider the following sequences: $abcde*abcd$ and $abcde*cdeab$, with the constraints ab , cde , and $abcd$. $*$ represents a sequence of arbitrary tokens. Consider a greedy algorithm that has processed the first five tokens ($abcde$), as shown below:

	abcde*
Interpretation 1:	ab cde
Interpretation 2:	abcd

Due to greediness, the algorithm must pick one of the two interpretations: the beginning of the string contains the constraints ab and cde , or it contains the constraint $abcd$. As the greedy algorithm does not backtrack, the end of the string cannot be taken into account. However, the correct interpretation depends on how the string ends. If the ending is $abcd$, the first interpretation was correct. On the other hand, if it is $cdeab$, the second interpretation was. Thus, all greedy algorithms fail to detect at least one of $abcde * abcd$ and $abcde * cdeab$.

The finite-state automaton-based approach suggested by Anderson et al. (2017) does not suffer from greediness, but is, as presented, infeasible in our use case: the number of hypotheses is 2^C , growing exponentially as the number of constraints C is increased.

2.3 Rare and Obscure Terms

Our method tackles what we have termed *obscure terms*. These are terms that are completely unpredictable to the model being decoded, as they are lexically surprising and have not been seen at training time. We argue that the challenge posed by obscure terms is fundamentally different from synonym selection in constrained decoding: not only are we willingly sampling along suboptimal paths, but we also sometimes have to choose tokens that directly contradict the analysis of the language model.

The nature-related lexica (Metsäkeskus 2022, 2023) we used during development and evaluation were full of these obscure terms. The Forest Centre terminologies contained such terms as “papintap-paja” (violet tanbark beetle, lit. “priest-slayer”) and “tukkimiehentäi” (large pine weevil, lit. “log-man’s louse”). These target terms are very unexpected and clearly unfamiliar to the model, as seen from the low term accuracy of the unconstrained translation model in our evaluation (see Section 4).

Obscure terms produced characteristic failures. As the model evaluates all positions for a constraint to be unlikely, it will often result in *misplacements* as well as *leaks* where a part of the structure of a source language constraint still made its way to the translation. See Table 2 for an example of both.

3 The Mitra Pipeline

We designed our constrained translation pipeline with three goals: 1) allowing a high number of disjunctive constraint alternatives to support highly agglutinative languages such as Finnish, 2) fix the problems caused by greediness in the previous constraint recognition algorithms, and 3) make constrained decoding a viable alternative even for rare terms not present in the training data of the neural network. The goals 1 and 2 are fulfilled by using a custom finite-state automaton-based constraint recognition algorithm, while goal 3 is fulfilled by a backtranslation substitution algorithm.

The full end-to-end pipeline contains the following components:

1. **Term Recognition and Constraint Generation.** A dependency parser is used to extract the noun, verb, and adjectival phrases contained in the input sentence. If any of the phrases is found in the lexicon, all supported inflected forms of the target-language term are generated and added as a constraint.

2. **Backtranslation Substitution.** Each of the target-language terms added as constraints is translated back to the source language using an NMT model trained on the same data as the model used to perform the actual translation. The input sentence is modified so that the recognized terms are replaced with the back-translations, inflected and capitalized similarly to the original terms.
3. **Constrained Beam Search.** A constrained beam search is performed to translate the input sentence to the target language.

We implement the pipeline for Finnish, Swedish, and English in all language directions. As our pipeline is agnostic to the NMT model itself, it can be used with any model as long as the appropriate language-specific modules have been implemented.

3.1 Term Recognition and Constraint Generation

Term recognition refers to scanning the input sentence and detecting all the phrases in the sentence that also appear in the terminology. We provide a phrase detection module for each of the supported languages. Each of the modules first performs the following high-level steps, although the specific methods are highly language-dependent and not within the scope of this paper.

1. **Dependency parsing.** We use the Stanford Stanza Python package (Qi et al. 2020) for English and Swedish, and the TranSmart dependency parser (Nykänen 1996) for Finnish.
2. **Phrase detection.** We iterate the dependency tree recursively and for each noun, adjective, and verb, we construct a list of noun phrases, adjectival phrases and verb phrases, respectively, as explained in the next step.
3. **Dependent selection.** For each noun, we iterate all combinations of its adjectival dependents (with the limit up to 6 adjectival dependents). We assemble a list of these combinations. For example, if the phrase is “young, strong cat”, we would generate the combinations “cat”, “young cat”, “strong cat”, and “young, strong cat”. For Finnish, we also return parts of compound nouns. For adjectives and verbs, we do not include any dependents, and simply return a one-item list with the adjective or verb itself.

		Recognized terms and backtranslations	
Original sentence	Hosat ovat käteviä työkaluja.	{hosa → fire swatter}	
Backtranslation	Paloswatterit ovat käteviä työkaluja.	{paloswatteri ← fire swatter}	
Translation	Fire swatters are handy tools.		
Original sentence	The characteristics of the live fuel type are defined mainly on the basis of the tree stand and ground vegetation.	{live fuel type palokasvustotyyppi}	→
Backtranslation	The characteristics of the type of fire growth are defined mainly on the basis of the tree stand and ground vegetation.	{type of fire growth palokasvustotyyppi}	←
Translation	Palokasvustotyyppin ominaisuudet määritellään pääasiassa puuston ja maakasvillisuuden perusteella.		

Table 3: Example of phrase detection and backtranslation substitution.

4. **Lemmatization.** We lemmatize each phrase returned by the previous step. For Finnish and Swedish, this includes taking into account adjective-noun agreement: Finnish nouns agree in case and number, Swedish nouns in determinateness, gender and number.

After phrase detection, we compare the list of phrases to the terminology and generate a list of constraints. For each term, the disjunctive constraint has multiple alternatives corresponding to the different inflectional forms of the term. Similarly to above, we perform dependency parsing for the target-language terms, and then inflect them taking into account adjective-noun agreement. For Finnish, we do not generate all the possible forms due to their high number, instead we have a list of the most common forms.

We use several open-source¹ and proprietary² language modules. Details of this step are not within the scope of this paper.

3.2 Backtranslation Substitution

After phrase detection, we perform backtranslation for all target language terms that have been added as constraints by using a reverse-language NMT model trained with the same dataset as the model used for translation proper. In our experiments, we

¹We use the Python packages *stanza*, *pyvoikko* (for lemmatization of Finnish compound words), *pyomorfi* (for inflecting Finnish verbs), *taivutin* (for inflecting Finnish nominals), *inflex* (for inflecting English). For Swedish, we use a proprietary statistical guesser based on the Saldo inflecto (<https://github.com/kielikone/saldo-inflecto>).

²Mostly TranSmart pipeline components (Nykänen 1996)

use the Opus-MT Tatoeba Challenge models (Tiedemann 2020) which include models for both language directions for most language pairs.

After producing the backtranslations, we verify that they improve the translation quality by using the NMT model to calculate scores for both (original term → target term) and (backtranslation → target term) pairs. If the pair with the backtranslation yields a better score, we replace the original term in the source language sentence with the backtranslation. We use dependency parsing and the language-specific modules detailed in the previous section to inflect the backtranslation in the same form as the original term. We also match the initial letter case.

An example of the backtranslation substitution process is given in Table 3.

3.3 Constrained Beam Search

We use a constrained beam search algorithm very similar to the one described by Post and Vilar (2018). The details of our algorithm are presented in Appendix A. The main differences to the previous algorithms are in our constraint recognition algorithm, detailed in the following section.

3.4 Constraint Recognition during Beam Search

For constraint recognition, we adapt a finite-state automata (FSA)-based approach similar to Anderson et al. (2017). We note that while the method used by Anderson et al. requires 2^C hypotheses, a different beam allocation strategy such as the one used in Grid Beam Search (Hokamp and Liu 2017) or Dynamic Beam Allocation (Post and Vilar 2018)

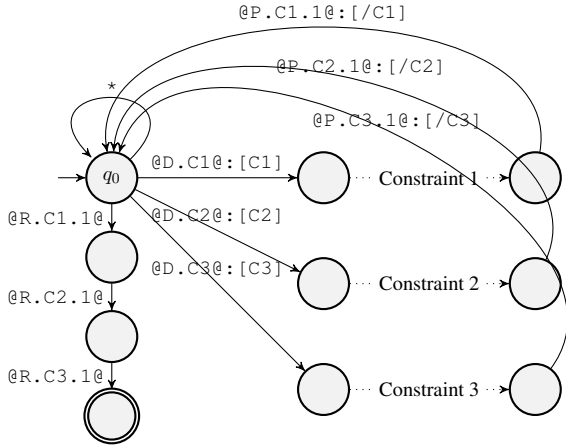


Figure 1: The structure of the finite state transducer that can recognize three constraints in any order. To reduce the size of the transducer, we use flag diacritics (Beesley and Karttunen 2003, chapter 8). $@D.f@$ is a flag diacritic that succeeds if f is undefined. $@P.f.1@$ is a flag diacritic that always succeeds and sets $f = 1$. $@R.f.1@$ is a flag diacritic that succeeds if $f = 1$. The nodes inside the constraints are not included; they would form a trie matching to all alternatives of that disjunctive constraint. The output of the transducer is the input, with symbols added for marking starts ($[C1]$) and ends ($[/C1]$) of constraints.

may also be used, which drops the beam size to either $O(C)$ or $O(1)$, respectively. In these scenarios, not every possible combination of constraints is stored in the beam: only the ones with the highest scores. We use a non-deterministic FSA that can track all possible interpretations at the same time: the two interpretations that would have been tracked in entirely different hypotheses in their solution can be tracked with a single hypothesis in our solution.

While the beam size can be limited to be linear with regard to the number of constraints, the memory constraints of the finite-state automaton cannot be. If the automaton is deterministic, its size is $O(2^C)$ in the worst-case scenario³. If the automaton is non-deterministic, its size will be $O(C)$, but the number of simultaneous states might be $O(2^{C \max |C_{i,j}|})$, i.e. in the worst case the number of states grows exponentially with regard to the number of constraints and their lengths. To mitigate this issue, we implement an optimization that removes most of the simultaneous states when we can safely determine that they recognize the same set of strings. This optimization is detailed below, after we have detailed the structure of the FSA.

³Since all possible combinations of recognized constraints (2^C) must be represented.

Figure 1 has a graph of a finite-state transducer that recognizes three constraints in any order. Since the number of orderings the constraints can be in is $C!$, we employ flag-diacritics (Beesley and Karttunen 2003, chapter 8) to reduce the number of nodes in the transducer. Flag diacritics behave like epsilon edges, but can only be followed if a variable, called a *flag*, is set to a specific value. In our transducer, each constraint has its own flag. The parts of the transducer matching to constraints are fenced with the $@D.f@$ diacritics that succeed only if the flag f is undefined, thus preventing the transducer recognizing any constraints more than one time. After a constraint is fully recognized, the $@P.f.1@$ diacritic is used to set the flag value to 1. The accepting node of the transducer is fenced with diacritics of type $@R.f.1@$ that require that the flag must be set to value 1.

The specific type of finite-state automata we use is the finite-state transducer, in which each edge can both consume an input symbol and output a symbol (Beesley and Karttunen 2003). By using a finite-state transducer instead of a regular finite-state machine, we detect the beginnings and ends of the constraints. The transducer outputs each input symbol, and additionally provides a start token (such as $[C1]$) and an end token (such as $[/C1]$) for each constraint. When we discuss “states” below, we refer to (q, m) pairs, in which q is a node in the finite-state machine, and m is the sequence of output symbols produced.

Since the transducer is non-deterministic, the number of parallel states can grow exponentially. For example, if there are C constraints and all of them are present in the input string, in the end there will be at least 2^C simultaneous states: one in which none of the constraints matched, one in which all of them matched, and all the possible combinations in between. To prevent this, we remove some of the states between each iteration based on the following condition: if two states S_1 and S_2 are both in q_0 (the initial state as in Figure 1), and the set of fulfilled constraints in the output of S_1 is a proper subset of fulfilled constraints in S_2 , the state S_1 is removed. The proof that this does not change the set of strings that are accepted by the FSA is included in Appendix B. This optimization makes the finite-state automaton computationally feasible on the real-world data we used.

The finite-state machines were implemented using the `kfst` Python package.

Configuration Name	Decoding Algorithm	Constraint Recognition	Backtranslations
Mitra-FB	CBS	NDFST	Yes
Mitra-TB	CBS	Trie	Yes
Mitra-F	CBS	NDFST	No
Mitra-T	CBS	Trie	No
Mixtral	Sampling ($T = 0.2$)	N/A	No
Poro	Sampling ($T = 0.2$)	N/A	No
Unconstrained	Greedy	N/A	No

Table 4: Evaluated configurations of the pipeline. CBS refers to our Constrained Beam Search algorithm as described in this paper. NDFST is our non-deterministic finite-state transducer. Trie refers to the constraint recognition algorithm inspired by Post and Vilar (2018); Hu et al. (2019) modified to support disjunctive constraints. In addition to our pipeline, we use Mixtral (Jiang et al. 2024) and Poro (SiloAI 2023), both large language models, and unconstrained machine translations. All methods apart from Mixtral and Poro use Opus-MT Tatoeba Challenge models for Finnish and English (Tiedemann 2020).

4 Evaluation

We use both automatic evaluation and human evaluation to measure the quality of the translations produced by our pipeline. The automatic methods include BLEU, chrF, TER, and COMET scores, as well as measuring the number of fulfilled constraints. In human evaluation, we asked a professional translator to evaluate all translated sentences and mark them either as OK, erroneous due to incorrectly applied constraints, or erroneous due to other cause.

Due to time constraints, we did not evaluate the Swedish translation even though we implemented it.

4.1 Evaluated Pipeline Configurations

For evaluation, we prepared four configurations of the Mitra pipelines: half of them use the NDFST-based constraint recognizer, and half of them a trie-based recognizer inspired by (Post and Vilar 2018; Hu et al. 2019) modified to support disjunctive constraints. Both of the two algorithms are evaluated with and without backtranslation substitution, with parameters $B = 1, k = 5, S = M = \infty$ and a 120 s timeout. In addition, we translated the sentences without constraints or backtranslations. In each case, we use the Opus-MT Tatoeba Challenge models for Finnish and English⁴ (Tiedemann 2020).

We also compare our methods to general-purpose language models Poro⁵ (SiloAI 2023) and Mixtral⁶ (Jiang et al. 2024) by embedding the con-

straint words into the prompt (cf. Ghazvininejad et al. 2023). For these models, we use the term recognition and constraint generation components of our pipeline, but do not apply backtranslations or enforce hard constraints. The prompts used are listed in Appendix C. We tried both zero-shot and 2-shot prompting and found that to achieve best results BLEU-wise, Mixtral needed to be prompted in a zero-shot manner and Poro in a 2-shot manner.

We ran the Opus-MT models on a Tesla T4 GPU, and the Poro and Mixtral models on an A100 80GB GPU.

We used greedy decoding with the unconstrained translation and temperature sampling with the LLMs, since these are the sampling methods most often used with these models. While beam search could have been used for both of these, and could have improved their performance, deciding the fair beam size would not have been trivial: the constrained translation has beam size $C + 1$, where C is the number of constraint tokens. If no constraints are used, this results in the beam size of 1, which corresponds to greedy decoding. Therefore, to simplify our experiment, we decided to not increase the beam size over this default unconstrained size of 1.

To save human evaluation resources, we performed the human evaluation for the model that received better BLEU scores, which is Mixtral for the Finnish–English translation direction, and Poro for the English–Finnish translation direction.

4.2 Evaluation Corpus

We conducted the evaluation on two vocabularies listed below. For both of them, we used the head words of the entries to construct the translation terminology and the definitions of the entries as the test sentences.

⁴<https://huggingface.co/Helsinki-NLP/opus-mt-tc-big-fi-en> and <https://huggingface.co/Helsinki-NLP/opus-mt-tc-big-en-fi>

⁵https://huggingface.co/LumiOpen/Poro-34B_700B_variant

⁶<https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1>

Configuration	Forest Fires (EN-FI)					Finnish Parliament (FI-EN)				
	BLEU	chrF	TER	COMET	Acc.	BLEU	chrF	TER	COMET	Acc.
Mitra-FB	20.96	62.31	69.62	0.90	100.00	35.50	63.77	53.57	0.85	98.34
Mitra-TB	19.70	59.35	72.10	0.89	94.83	35.70	64.33	53.48	0.84	100.00
Mitra-F	19.48	61.87	70.09	0.89	100.00	35.21	63.67	54.45	0.85	98.34
Mitra-T	19.13	59.19	75.06	0.89	94.83	35.40	64.12	54.61	0.84	100.00
Mixtral	7.71	46.49	92.43	0.79	84.48	32.99	62.85	55.71	0.84	82.82
Poro	16.01	55.95	79.20	0.88	83.62	21.13	55.38	70.71	0.84	50.27
Unconstrained	16.45	54.55	76.95	0.86	22.41	34.48	62.08	53.52	0.86	49.46

Table 5: The results for automatic evaluation of the configurations described in Table 4. “Acc.” refers to the number of fulfilled constraints.

Configuration	Forest Fires (EN-FI)			Finnish Parliament (FI-EN)		
	OK %	Constr. error %	Other error %	OK %	Constr. error %	Other error %
Mitra-FB	60.00	2.35	37.65	60.56	10.83	28.61
Mitra-TB	55.29	4.70	40.00	57.78	13.61	28.61
Mitra-F	55.29	9.41	35.29	56.39	15.00	28.61
Mitra-T	52.94	9.41	37.65	54.17	17.78	28.06
Mixtral				66.11	0.56	33.33
Poro	52.94	2.35	44.71			
Unconstrained	32.94	0.00	67.06	67.78	0.83	31.39

Table 6: The results for human evaluation of the configurations described in Table 4. All of the sentences were categorized to the three categories “OK”, “Erroneous due to incorrectly applied constraints” (Constr. err), and “Erroneous due to other error”. The timed out sentences are counted towards other errors.

1. **Forest Fire Vocabulary** by the Finnish Forest Centre (Metsäkeskus 2022), consisting of 85 Finnish/English word pairs and definitions. For this vocabulary, we translated the definitions in the English–Finnish direction.
2. **Finnish Parliament Vocabulary** by the Finnish Parliament (Eduskunta 2008), consisting of 360 Finnish/English word pairs and definitions. We used only 358 of these since two contained special characters for which the pre-processing pipeline failed. For this vocabulary, we translated the definitions in the opposite direction: Finnish–English.

We release all term pairs and test sentences openly⁷.

4.3 Evaluation Methods

We performed both automatic and human evaluation. For automatic evaluation, we calculated BLEU, chrF, and TER scores for the sentences using the `sacrebleu` Python library (Post 2018), and the COMET score⁸ (Rei et al. 2022) using the

⁷<https://github.com/kielikone/mitra-eval-results>

⁸The `wmt22-comet-da` model

`evaluate` Python library. Furthermore, we used the term recognition component of our pipeline to analyze the number of constraints fulfilled in the output sentences. This is similar to the lemmatized term exact match accuracy (Bergmanis and Pinnis 2021) and exact match accuracy (Alam et al. 2021a), although we do not need to specifically lemmatize the words as our disjunctive constraints include the inflected forms.

For manual evaluation, we generated a spreadsheet that contained one input sentence on each row, the reference translation, and the outputs of each of our tested configurations in random order. We asked a professional translator to evaluate each configuration and mark it either as correct, erroneous due to incorrectly applied constraints (while still present), or erroneous due to other cause (incl. missing constraint) (cf. Bergmanis and Pinnis 2021, ⁹). We then calculate percentages of these three categories for each of the evaluated configurations.

⁹Our evaluation differs from that of Bergmanis and Pinnis (2021): They had categories “wrong lexeme” and “wrong inflectional form”. We measure the presence of the constraint lexeme automatically, and are more interested in errors caused by incorrect placement of the constraint than those of wrong inflectional form, since the former errors are common in our tests and much more critical.

For the Parliament dataset, the NDFST-based methods (Mitra-FB and Mitra-F) timed out for three sentences, i.e. the beam search never reached the end condition. Similarly, for the Forest Fire Dataset, the trie-based method timed out for one sentence. Those sentences are evaluated as empty strings in the automatic evaluation and left out of the manual evaluation. See Appendix D for an analysis of them.

4.4 Results

The results of the automatic evaluation are presented in Table 5 and the human evaluation in Table 6.

For both datasets, the usage of constraints improved the BLEU, chrF, and TER scores when compared to the unconstrained translations and the general-purpose language model outputs. For the Forest Fire dataset, the BLEU of the unconstrained translations was 16.45, while the BLEU of Mitra-FB was 20.96. For the Parliament dataset, the unconstrained BLEU improved from 34.48 to 35.50 respectively. The COMET score improved from 0.86 to 0.90 with the Forest Fire dataset, while it decreased insignificantly from 0.86 to 0.85 with the Parliament dataset. All automatic evaluation scores for the different Mitra configurations were too near each other to be significant.

Similarly, Mitra-F and Mitra-FB raised the number of fulfilled constraints to 100% from 22.41% for the Forest Fire dataset, with the trie-based methods Mitra-T and Mitra-TB timing out with one sentence, causing the percentage to drop to 94.83%. For the Parliament dataset, the NDFST-based methods timed out for three sentences, causing the fulfilled constraint percentage to reach only 98.32%, while the unconstrained translations reached 49.46%.

The general-purpose language models achieved BLEU scores comparable to the unconstrained translation with the exception of Mixtral in the English–Finnish direction that produced translations of unusably low quality. Similarly, both Poro and Mixtral fulfilled ca. 82–84% of the constraints with the exception of Poro in the Finnish–English direction. Since Mixtral is arguably better when English is the target language, and Poro when Finnish is the target language, we did not conduct human evaluation for both models, choosing instead the model that performed better for the evaluated language direction.

Configuration	Mean time	
	Failed excl.	Failed = 120 s
Mitra-FB	2.11 s	2.11 s
Mitra-TB	1.68 s	3.07 s
Mitra-F	2.36 s	2.36 s
Mitra-T	1.44 s	2.83 s
Unconstrained	0.16 s	0.16 s

Table 7: Mean translation times for the Forest Fire datasets. Mixtral and Poro times are not included since they were ran on a different hardware. There was one sentence that timed out with the trie-based configurations. In the first column, that sentence was removed. In the second column, that sentence was given the value equal to the timeout we used, 120 s.

In human evaluation, the NDFST-based approach and backtranslations achieve significantly better results than the Trie and non-backtranslated configurations. For the Forest Fire dataset, they together raise the number of “OK” translations from 52.94% to 60.00%. Similarly, for the Parliament dataset, the number rose from 54.17% to 60.56%. At the same time the number of errors caused by incorrectly applied constraints decreased from 9.41% to 2.35%, and from 17.78% to 10.83%, respectively. For the Forest Fire dataset, the usage of constraints increased the quality when compared to unconstrained translations, while for the Parliament dataset, the unconstrained sentences were evaluated to have higher quality.

For both datasets, the number of “other errors” was considerable. For the Forest Fire dataset, our evaluator noted that the Forest Fire vocabulary did not contain all of the special jargon used in the sentences. Thus, had the vocabulary been more comprehensive, the translation quality could have been better.

4.5 Time performance

The mean translation times for the Forest Fire dataset sentences are presented in Table 7. For constrained translations, the time was measured for the full pipeline, including preprocessing and dependency parsing. The unconstrained translation times are significantly lower than the constrained times, but the time does not include any preprocessing.

Of the sentences which both algorithms were able to translate, the NDFST-based configurations were slower than the trie-based configurations. However, as the trie-based algorithm failed to translate one of the sentences (due to its greediness, it was unable to place the constraints to the sentence correctly, which led to the NMT model considering

all hypotheses improbable and never finishing¹⁰). As the translation timed out, one might argue that the algorithm should be penalized for this by counting the sentence using the timeout as the time it took to “produce” the empty translation. We have reported both numbers in different columns.

We ran Mixtral, Poro, and the Parliament dataset evaluations on different hardware and software environments, so we cannot present comparable numbers for them. We used no batching of multiple sentences.

5 Discussion and Conclusions

In this paper, we have presented *Mitra*, a pipeline for terminologically-constrained machine translation that improves on the previous “hard” constraint methods with a finite state automaton-based constraint recognition algorithm and a backtranslation substitution step. When compared to the trie-based method based on the previously suggested algorithms (Post and Vilar 2018; Hu et al. 2019) without backtranslations, our method significantly increases the quality in human evaluation.

We argue that our method fulfills the three goals we began with: allowing disjunctive constraints, solving the problems of greedy constraint recognition algorithms, and improving quality on “rare and obscure” terms. The finite state automata-based algorithm allows any number of alternatives in the disjunctive constraints, and does not suffer from the greediness of previous algorithms. Furthermore, when combined with backtranslations, the number of errors caused by incorrectly applied constraints drops significantly on both of the evaluation datasets.

While constrained translation significantly improved quality on the Forest Fire dataset, it unexpectedly decreased quality on the Parliament dataset in human evaluation. We believe this discrepancy is due to two main factors. Firstly, the Parliament dataset’s vocabulary is not well-suited for constrained generation (see Appendix D for examples), and secondly, the Opus-MT models we used are more familiar with the subject matter, using 49.46% of the constraint terms even when unconstrained, leaving less room for improvement. The Forest Fire dataset, on the contrary, is very unfamiliar to the model, as only 32.94% of the sentences were translated acceptably, and only 22.41% of the

desired target terms were used when unconstrained. The effectiveness of constrained translation thus depends on the quality of the constraint vocabulary and the topic of the texts translated.

The major downside of “hard” constraints is their increased time requirement: the translation times were more than ten times larger on three of the four configurations evaluated (see Table 7). Although we did not optimize the evaluation by using batching or other methods such as those recommended by Hu et al. (2019), it is clear that constrained translation is slower in any case. Of our configurations, those that use the NDFST algorithm are slower than those that use tries. However, this is not as major a problem as one might initially think, as in most cases the trie-based solution yields the same result as the NDFST solution. The problem of greediness is only present when two constraints share tokens – if no constraints overlap in this way, there is no ambiguity. Thus, when translating a longer text, the trie-based approach can be used instead of the NDFST solution for most input sentences, making the translation of the whole text nearly as fast as when translated completely with the trie method.

Of the general-purpose large language models evaluated, Mixtral performed very well on the Parliament dataset, producing higher-quality results than constrained decoding methods. Similarly, Poro provides results comparable to the *Mitra-T* configuration in the human evaluation. Although the percentage of fulfilled constraints is lower than that of constrained translation, ca. 82–84%, the number of “constraint errors” is also low, implying that at least some of the missing constraints can be explained by the models providing a satisfying translation using a synonym or other acceptable construct that does not match the constraint when they have trouble fitting the constraint word into the sentence. Thus, the error mode of LLMs might be considered better than that of constrained decoding-based translation.

Since our pipeline is model-agnostic, it can be used with any NMT model or even with a general-purpose language model. Similarly, many of the “soft” constraint methods can also be combined with our method. We believe that our future research should focus on evaluating these combinations. As even a considerably basic soft method such as backtranslations can improve the translation quality significantly, our hypothesis is that more complex soft methods (such as Bergmanis and Pinnis 2021) can improve it even further.

¹⁰See Table 1 for a simplified version of the sentence and Appendix D for a full analysis of the failed sentences

Acknowledgments

Both IH and TF are funded by Kielikone Oy. IH is additionally funded by the Doctoral Program in Computer Science at the University of Helsinki. We would like to thank Piia Saresoja for evaluating the translated sentences and Kaarina Hyvönen for lending time towards polishing this text.

References

- Alam, Md Mahfuz Ibn, Antonios Anastasopoulos, Laurent Besacier, James Cross, Matthias Gallé, Philipp Koehn, and Vassilina Nikoulina. 2021a. On the evaluation of machine translation for terminology consistency. *Preprint*, arXiv:2106.11891.
- Alam, Md Mahfuz Ibn, Ivana Kvapilíková, Antonios Anastasopoulos, Laurent Besacier, Georgiana Dinu, Marcello Federico, Matthias Gallé, Kweonwoo Jung, Philipp Koehn, and Vassilina Nikoulina. 2021b. Findings of the wmt shared task on machine translation using terminologies. In *Proceedings of the Sixth Conference on Machine Translation*, pages 652–663.
- Anderson, Peter, Basura Fernando, Mark Johnson, and Stephen Gould. 2017. Guided open vocabulary image captioning with constrained beam search. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 936–945, Copenhagen, Denmark. Association for Computational Linguistics.
- Arnola, Harri. 1996. Kielikone Finnish-English MT system “TranSmart” in practical use. In *Proceedings of Translating and the Computer 18*.
- Beesley, Kenneth R and Lauri Karttunen. 2003. Finite-state morphology: Xerox tools and techniques. *CSLI, Stanford*, pages 359–375.
- Bergmanis, Toms and Mārcis Pinnis. 2021. Facilitating terminology translation with target lemma annotations. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3105–3111, Online. Association for Computational Linguistics.
- Eduskunta. 2008. Valtioneuvoston termipankki Valter: Eduskuntasanasto.
- Ghazvininejad, Marjan, Hila Gonen, and Luke Zettlemoyer. 2023. Dictionary-based phrase-level prompting of large language models for machine translation. *Preprint*, arXiv:2302.07856.
- Hasler, Eva, Adrià de Gispert, Gonzalo Iglesias, and Bill Byrne. 2018. Neural machine translation decoding with terminology constraints. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 506–512, New Orleans, Louisiana. Association for Computational Linguistics.
- Hokamp, Chris and Qun Liu. 2017. Lexically constrained decoding for sequence generation using grid beam search. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1535–1546, Vancouver, Canada. Association for Computational Linguistics.
- Hu, J Edward, Huda Khayrallah, Ryan Culkin, Patrick Xia, Tongfei Chen, Matt Post, and Benjamin Van Durme. 2019. Improved lexically constrained decoding for translation and monolingual rewriting. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 839–850.
- Jiang, Albert Q., Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Léo Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. Mixtral of experts. *Preprint*, arXiv:2401.04088.
- Koehn, Philipp. 2009. *Statistical machine translation*. Cambridge University Press.
- Koehn, Philipp, Franz J. Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 127–133.
- Metsäkeskus. 2022. Tuli metsässä -sanasto. Received: 2023-06-20.

Metsäkeskus. 2023. Tuhonaiheuttajat -sanasto. In prep. Received: 2023-06-02.

Nykänen, Asko. 1996. Design and Implementation of an Environment for Parsing Finnish Sentences.

Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Post, Matt and David Vilar. 2018. Fast lexically constrained decoding with dynamic beam allocation for neural machine translation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1314–1324.

Qi, Peng, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.

Rei, Ricardo, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Semenov, Kirill, Vilém Zouhar, Tom Kocmi, Dongdong Zhang, Wangchunshu Zhou, and Yuchen Eleanor Jiang. 2023. Updated findings of the wmt 2023 shared task on machine translation with terminologies. In *Proceedings of the Eight Conference on Machine Translation (WMT)*. Association for Computational Linguistics. The updated version available at https://wmt-terminology-task.github.io/upd_wmt_terminology_2023.pdf.

SiloAI. 2023. Poro - a family of open models that bring European languages to the frontier. <https://www.silo.ai/blog/poro-a-family-of-open-models-that-bring-european-languages-to-the-frontier>. Accessed: 2024-02-20.

Tiedemann, Jörg. 2020. The Tatoeba Translation Challenge – Realistic data sets for low resource

and multilingual MT. In *Proceedings of the Fifth Conference on Machine Translation*, pages 1174–1182, Online. Association for Computational Linguistics.

A Algorithms

We have described our beam search algorithm in detail here for increased reproducibility. Algorithm 1 contains the main loop of the beam search. Algorithm 2 contains the candidate selection algorithm. Finally, Algorithm 3 details the beam allocation algorithm.

Since our algorithm is model-agnostic, we have abstracted the calls to the NMT model by referring to the probability function P and “top- k sampling”. In real implementations, these probabilities would come from the decoder of the NMT model or an LLM model.

```

input :  $M$  maximum length of the output in tokens
          $C$  number of constraint tokens
          $B$  hypothesis bank size
          $S$  beam size
          $V$  the vocabulary
output : best hypothesis

hypotheses  $\leftarrow$  [[start token]]
cutoff  $\leftarrow$  0
best hypothesis  $\leftarrow$  null
for  $M$  times do
  /* Calculating the new hypotheses */
  candidates  $\leftarrow$  GetCandidates(hypotheses,  $V$ ,  $k$ )
  hypotheses  $\leftarrow$  Allocate(candidates,  $C$ ,  $B$ ,  $S$ )
  /* Updating cutoff */
  foreach finished hypothesis  $h$  in hypotheses do
    if  $P(h) >$  cutoff then
      cutoff  $\leftarrow$   $P(h)$ 
      best hypothesis  $\leftarrow$   $h$ 
    end
  end
  /* Pruning hypotheses */
  foreach hypothesis  $h$  in hypotheses do
    if  $P(h) \leq$  cutoff then
      remove  $h$  from hypotheses
    end
  end
  /* Early stopping */
  if |hypotheses| = 0 then
    return best hypothesis
  end
end

```

Algorithm 1: The main loop of the beam search.

B Proofs

The “NDFST” in these proofs refers to a non-deterministic finite-state transducer that has the structure described in this paper (see Figure 1). q_0 is the initial state of the NDFST and q_a is the sole accepting state.

Definition. $C(m)$ is the set of constraint end tokens in the output symbol list m . Since the flag for a

```

Function GetCandidates (hypotheses,  $V$ ,  $k$ ) is
  input :hypotheses from the previous iteration
           $V$  the vocabulary
           $k$  the parameter for top- $k$  sampling
  output :set of candidate hypotheses
  candidates  $\leftarrow \emptyset$ 
  /* Constraint continuations from the NDFST
  */
  foreach  $h$  in hypotheses do
    if  $h$  is not finished then
      foreach token  $t$  that would advance the current
        NDFST state do
        append  $h + t$  to candidates
      end
    end
  end
  /* Unconstrained candidates
  */
  foreach  $h$  in hypotheses do
    append top- $k$  continuations of  $h$  to candidates
  end
  return candidates
end

```

Algorithm 2: The candidate selection algorithm. This is a simplified version of the algorithm described in (Post and Vilar 2018, section 3.1), combining their steps 1 and 3.

```

Function Allocate (candidates,  $C$ ,  $B$ ,  $S$ ) is
  input :the list of candidates
           $C$  number of constraint tokens
           $B$  target hypothesis bank size
           $S$  maximum beam size
  output :a new list of hypotheses
  if  $B \cdot (C + 1) > S$  then
    return Allocate2 (candidates,  $C$ ,  $\lfloor \frac{S-1}{C} \rfloor + 1$ ,  $S$ )
  else
    return Allocate2 (candidates,  $C$ ,  $B$ ,  $B \cdot (C + 1)$ )
  end
end
Function Allocate2 (candidates,  $C$ ,  $B'$ ,  $S'$ ) is
  input :the list of candidates
           $C$  number of constraint tokens
           $B'$  actual hypothesis bank size
           $S'$  actual beam size
  output :a new list of hypotheses
  /* Allocate hypotheses to banks they
  belong based on their number of
  fulfilled constraint tokens
  */
  hypotheses  $\leftarrow []$ 
  foreach  $i$  in  $C, \dots, 0$  do
    bank size  $\leftarrow 0$ 
    foreach candidate  $c$  from most probable to least probable
      do
        if bank size  $< B'$  then
          if number of constraint tokens in  $c \geq i$  then
            append  $c$  to hypotheses
            remove  $c$  from candidates
            bank size  $\leftarrow$  bank size + 1
          end
        end
      end
    end
  end
  /* Fill underfilled banks with most
  probable candidates
  */
  foreach  $i$  in  $C, \dots, 0$  do
    foreach candidate  $c$  from most probable to least probable
      do
        if |hypotheses| =  $S'$  then
          return hypotheses
        end
        if number of constraint tokens in  $c \geq i$  then
          append  $c$  to hypotheses
          remove  $c$  from candidates
        end
      end
    end
  end
  return hypotheses
end

```

Algorithm 3: The beam allocation algorithm.

constraint is set in the same transition that generates the constraint end token, $C(m)$ also corresponds to the set of constraints that have their flag set. C is the set of all constraint end tokens possible.

Definition. $S(q, m)$ is the set of strings that the NDFST accepts from the initial state q and the initial output symbol list m .

Theorem. Given a string s partitioned into two parts $s_1 s_2$, so that the NDFST has consumed s_1 but not s_2 , and NDFST states (q_0, m_1) and (q_0, m_2) , so that $C(m_1) \subsetneq C(m_2)$ and $s_2 \in S(q_0, m_1)$, then $s_2 \in S(q_0, m_2)$.

Proof. The only accepting node in the NDFST is fenced with flag diacritic symbols, each corresponding to a different constraint. Thus, the accepting state can only be reached if all the flags are set, that is, $C(m) = C$.

Since $s_2 \in S(q_0, m_1)$, there must be a path $q_0 \rightarrow q_{i_1} \rightarrow \dots \rightarrow q_{i_n} \rightarrow q_a$ accepted by the NDFST if (q_0, m_1) is used as an initial state.

As $C(m_1) \subsetneq C(m_2)$, the same path cannot be accepted when (q_0, m_2) is used as the initial state, as for each $c \in C(m_2) \setminus C(m_1)$, the path contains a negative flag diacritic check that prevents the path from being accepted, as the flag for c is already set since $c \in C(m_2)$.

We construct a new path that is accepted when (q_0, m_2) is used as the initial state. For each $c \in C(m_2) \setminus C(m_1)$, we modify the path by replacing each transition beginning from $@D.Cc@ : [Cc]$ and ending to $@P.Cc.1@ : [/Cc]$ with $q_0 \rightarrow q_0$. As all of these c were already present in $C(m_2)$, this modification only removes duplicate positive flag diacritic sets. Since all the flags are set in the modified path, it is accepted by the NDFST. QED.

C LLM Prompts

C.1 Finnish–English

Please translate the following sentence using this vocabulary. Respond using JSON output such as {"translation": "This is the translation"}.

Vocabulary: joki = river; virtaava vesi = flowing water; valuma-alue = catchment area

Sentence: Vesilaissa joella tarkoitetaan virtaavan veden vesistöä, jonka valuma-alue on vähintään sata neliökilometriä.

C.2 English–Finnish

Käännä lause suomeksi annetulla sanastolla. Vastaa JSON-muodossa, esim. {"käännös": "Tämä on käännös"}.
 Sanasto: octopodes = mustekalat; extant = elävä; subclass = alaluokka; cephalopod = pääjalkainen; nautilus = helmivene
 Lause: Octopodes are one of the two extant subclasses of the cephalopods. It is also called two-gilled cephalopods. The other subclass is the nautilus or four-gilled cephalopods.

D Failed Translations

Several test sentences timed out during the evaluation. This section contains these sentences as well as an analysis of the cause of the error.

To save space, when we list the constraints, we only list the target lemmas, although in reality we give all inflections generated by our phrase inflector module as a disjunctive constraint. The disjunctive constraint also includes differently capitalized versions of the target term, although we list all target terms in lower case here.

D.1 Forest Fires

This evaluation corpus had only one failed sentence, which failed for both trie-based configurations, but not for the NDFST-based configurations.

D.1.1 Sentence 46

Constraints are underlined.

Source sentence: Forest fire regime describes the role of fire in a given area over a given time period and sums up almost all variables related to forest fires: forest fire effects and their influencing factors, forest fire frequency, forest fire severity, forest fire intensity, the size of fire, the time of and reason for ignition, regularity, variation, etc.

Reference translation: Metsäpaloregiimi kuvaa tulen roolia tietyssä aikana tietyllä alueella ja summaa yhteen lähes kaikki metsäpaloihin liittyvät suuret: metsäpalojen vaikutukset ja niihin vaikuttavat tekijät, metsäpalojen toistuvuus, metsäpalon vaikuttavuus, metsäpalon voimakkuus, koko, syttymisajankohta, syttymissy, säännöllisyys, vaihtelu, jne.

Constraints: {forest fire regime → metsäpaloregiimi}, {forest fire → metsäpalo}, {forest fire → metsäpalo}, {forest fire frequency → metsäpalojen toistuvuus}, {forest fire severity

→ metsäpalon vaikuttavuus}, {forest fire intensity → metsäpalon voimakkuus}

The timeout of this sentence is caused by two factors: the greediness of the constraint recognition algorithms and the large number of constraints it has. See Table 1 for an example of a simplified version of this sentence that does not timeout since it has less constraints.

The core issue is that after the tokens encoding “Metsäpalo” have been generated, the trie-based algorithm marks the constraint {forest fire → metsäpalo} fulfilled. After this, the beam search generates the tokens encoding “regiimi”, but they are not recognized to be a part of the constraint, since the progress of all other constraints was reset when one of the possible constraints tracked simultaneously was marked fulfilled. This means that the constraint {forest fire regime → metsäpaloregiimi} is yet unfulfilled, and the algorithm tries to place it later in the sentence. However, the decoder language model (correctly) considers all those other places to be improbable, and thus the end condition of the beam search is never reached within the time limit.

D.2 Finnish Parliament

This evaluation corpus had two sentences for which the preprocessing pipeline (i.e. morphological analysis and generation) failed due to a bug that we had not time to correct. We left these sentences out of the evaluation. Of the remaining sentences, the NDFST-based configurations failed for three sentences.

These sentences fail due to two main reasons: they contain too many constraints, and the target terms included in the vocabulary are poorly suited to be used as constraints resulting in unnatural translations considered improbable by the decoder language model. Note that while this section only analyses the timed out sentences, the poor suitability of the vocabulary applies also to those sentences that did not time out and is one cause to the poor scores received by the system in human evaluation.

D.2.1 Sentence 150

Constraints are underlined. This sentence was actually a text with two sentences, the second of which was more problematic.

Source sentence: Valtiopäiväasiakirjat julkaistaan painettuina Valtiopäiväasiakirjat-sarjassa sekä

nykyään myös eduskunnan sivustolla Internetissä. Valtiopäiväasiakirjat-sarjassa julkaistaan mm. eduskunta-aloitteet, eduskunnan täysistuntojen pöytäkirjat ja niiden ruotsinkieliset lyhennelmät, hallituksen esitykset, valtioneuvoston kirjelmät, tiedonannot ja selonteot, valiokuntien mietinnöt ja lausunnot, eduskunnan vastaukset ja kirjelmät, välikysymykset sekä kirjalliset kysymykset vastauksineen.

Reference translation: Parliamentary documents are published in print form in the series ‘Valtiopäiväasiakirjat’ and in recent years have been published on Parliament’s web pages as well. The series contains, among other documents, parliamentary motions; records of plenary sessions of Parliament and their Swedish summaries; government proposals, communications, statements and reports; committee reports and statements; parliamentary replies and communications; interpellations; and written questions and the replies to them.

Constraints for the first sentence:

{valtiopäiväasiakirja → parliamentary document}, {eduskunta → parliament, finnish parliament, eduskunta}

Constraints for the second sentence: {eduskunta-aloite → parliamentary motion, member of parliament’s motion, member’s motion}, {eduskunta → parliament, finnish parliament, eduskunta}, {täysistunto → plenary session}, {hallituksen esitys → government proposal}, {valtioneuvoston kirjelmä → government communication}, {valiokunta → committee}, {mietintö → report of the committee, committee report}, {lausunto → statement of the committee, committee statement}, {eduskunnan vastaus → parliamentary reply}, {välikysymys → motion of censure, interpellation}, {kirjallinen kysymys → written question}

The primary reason for the timeout of this text on the NDFST-based configurations is the large number of constraints. The NDFST-based algorithm has exponential time complexity in the worst-case scenario.

This text is also problematic due to the poor suitability of the vocabulary for constrained translation. For example, the phrase “valiokuntien mietinnöt ja lausunnot” (“the reports and statements of the committees”) generates the constraints {valiokunta → committee}, {mietintö → report of the committee, committee report}, and

{lausunto → statement of the committee, committee statement}, all of which include the word “committee” (valiokunta). Thus the constraints force the beam search hypotheses to contain translations like “the reports of the committee and the statements of the committee of the committees”, which the decoder language model obviously considers improbable. In some hypotheses the extraneous “committee” words also appear in completely different (and wrong) places in the translation, causing hallucinations.

D.2.2 Sentence 232

Constraints are underlined.

Source sentence: Eduskunnan tilintarkastajat antavat eduskunnalle kaksi tilintarkastuskertomusta: 1) eduskunnan tilintarkastajien tilintarkastuskertomuksen eduskunnan tilinpäätöksestä, toimintakertomuksesta ja kirjanpidosta sekä hallinnosta ja 2) eduskunnan tilintarkastajien tilintarkastuskertomuksen Valtiontalouden tarkastusviraston tilinpäätöksestä, toimintakertomuksesta ja kirjanpidosta sekä hallinnosta.

Reference translation: The parliamentary auditors submit two reports to Parliament: 1) a report on the financial statements, annual report and accounting, and administration of Parliament; and 2) a report on the financial statements, annual report and accounting, and administration of the National Audit Office of Finland.

Constraints: {eduskunta → parliament, finnish parliament, eduskunta}, {eduskunta → parliament, finnish parliament, eduskunta}, {tilintarkastuskertomus → parliamentary auditors’ report, report of the auditors of parliament}, {eduskunta → parliament, finnish parliament, eduskunta}, {tilintarkastuskertomus → parliamentary auditors’ report, report of the auditors of parliament}, {eduskunta → parliament, finnish parliament, eduskunta}, {eduskunta → parliament, finnish parliament, eduskunta}, {tilintarkastuskertomus → parliamentary auditors’ report, report of the auditors of parliament}, {valtiontalouden tarkastusvirasto → national audit office of finland}

As with the previous sentence, this sentence contains a large number of constraints. However, un-

like in the previous case, here most of the constraints are identical. In fact, it only contains three unique constraints. However, all constraints, even if duplicate, will get separate paths in the finite state machine. The algorithm could be easily optimized by adding counters instead of binary flags (cf. Hu et al. 2019). However, the FST library we used did not support them and we did not have time to implement them.

Further issues are caused by the fact that the correct translation of “eduskunnan tilintarkastaja” is “parliamentary auditor”. However, since it is not included in the vocabulary, the only constraint added is {eduskunta → parliament, finnish parliament, eduskunta}. The “parliament” added as a constraint clashes with the correct word “parliamentary” (in our case, both are single tokens, so they don’t share subword tokens), causing hallucinations in some hypotheses as “parliament” is inserted into a wrong place in the sentence, although the best hypothesis in this case includes arguably passable “the auditors of Parliament”.

Again, as in the previous case, the translations given in the vocabulary are unsuitable for constrained translation. The phrase “eduskunnan tilintarkastajien tilintarkastuskertomuksen” generates the constraints {eduskunta → parliament, finnish parliament, eduskunta} and {tilintarkastuskertomus → parliamentary auditors’ report, report of the auditors of parliament}, leading to unnatural translations such as “the parliamentary auditors’ report of the Parliament”. In fact, the literal translation of “eduskunnan tilintarkastajien tilintarkastuskertomus” is “parliamentary auditors’ report”, i.e. the translation given to the last word in the phrase is the translation of the whole phrase. To be suitable for constrained generation, the vocabulary should contain only one-to-one equivalent translations.

D.2.3 Sentence 321

Constraints are underlined.

Source sentence: Jos kansanedustaja kesken vaalikauden kuolee, hänelle myönnetään vapautus tai hänet erotetaan kokonaan edustajantoimestaan tai hän siirtyy Euroopan parlamentin jäseneksi, hänen tilalleen eduskuntaan tulee varaedustaja joko vaalikauden loppuun saakka tai määrääjäksi.

Reference translation: If a Member of Parliament dies during the electoral term, is granted a release from office, is dismissed from office, or is elected

to the European Parliament, he or she is replaced in Parliament for the remainder of the electoral term or for a specific period of time by a replacement Member.

Constraints: {kansanedustaja → member of parliament, representative, mp}, {vaalikausi → term of parliament, parliamentary term, electoral term}, {edustajantoimi → office of representative, mp’s responsibilities, member’s responsibilities}, {parlamentti → parliament, legislature}, {eduskunta → parliament, finnish parliament, eduskunta}, {varaedustaja → replacement member of parliament, alternate member of parliament, deputy member of parliament}, {vaalikausi → term of parliament, parliamentary term, electoral term},

Again, this sentence has many constraints. Also, like with the previous two cases in the Finnish Parliament dataset, the vocabulary used to generate the constraints is unsuitable for constrained translation. Since the translation of “edustajantoimi” must be either “office of representative”, “mp’s responsibilities” or “member’s responsibilities”, the phrase “Jos kansanedustaja [...] erotetaan kokonaan edustajantoimestaan” (“If a Member of Parliament [...] is dismissed from office”) must be translated clumsily as “If a Member of Parliament [...] is removed from the office of Representative”, where the word “office of Representative” is unnecessarily used instead of just “office”.

Similarly, the translations of the term “varaedustaja” (lit. “replacement representative”) are “replacement member of parliament”, “alternate member of parliament”, or “deputy member of parliament”, all of which contain unnecessarily the word “parliament”. Thus, the phrase “hänen tilalleen eduskuntaan tulee varaedustaja” (“he or she is replaced in Parliament [...] by a replacement Member”) must be translated with “he or she is replaced in Parliament by a replacement Member of Parliament”, duplicating the word “Parliament”.

If the constraints force the translator to generate unnatural text, the decoder will again give low scores to all hypotheses, thus making it difficult to reach the end condition in time.

Bootstrapping Pre-trained Word Embedding Models for Sign Language Gloss Translation

Euan McGill
Universitat Pompeu Fabra
Barcelona, Spain
euan.mcgill@upf.edu

Luis Chiruzzo
Universidad de la República
Montevideo, Uruguay
luischir@fing.edu.uy

Horacio Saggion
Universitat Pompeu Fabra
Barcelona, Spain
horacio.saggion@upf.edu

Abstract

This paper explores a novel method to modify existing pre-trained word embedding models of spoken languages for Sign Language glosses. These newly-generated embeddings are described, visualised, and then used in the encoder and/or decoder of models for the Text2Gloss and Gloss2Text task of machine translation. In two translation settings (one including data augmentation-based pre-training and a baseline), we find that bootstrapped word embeddings for glosses improve translation across four Signed/spoken language pairs. Many improvements are statistically significant, including those where the bootstrapped gloss embedding models are used.

Languages included: American Sign Language, Finnish Sign Language, Spanish Sign Language, Sign Language of The Netherlands.

1 Introduction

There has been a surge in research interest on Sign Language machine translation (SLMT) in recent years, but the data scarcity problem (De Sisto et al., 2022) and lack of standardised annotated data (Cormier et al., 2016) remain substantial obstacles to overcome.

At the heart of the labelling problem is the fact that although writing and transcription systems exist for SLs (Grushkin, 2017), none are used day-to-day by signers. Glosses are a semantic labelling

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

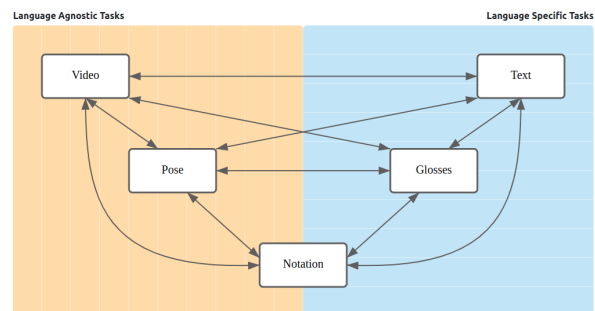


Figure 1: Intermediate, or subtasks of SLMT (Moryossef and Goldberg, 2021). This work focuses on translation between text and glosses.

tool for signs. They typically use lexemes from the *ambient* spoken language of the hearing community where the SL is used in order to convey the semantic sense of a given sign. However, glosses cannot be considered an orthographic system for SLs as they often differ between datasets, are not used by signers to write their languages (Müller et al., 2023), and may not include linguistic phenomena which are crucial to understand an utterance (Yin and Read, 2020).

SLMT is inherently multimodal (Bragg et al., 2019), and it is helpful to conceptualise it as a constellation of sub-tasks at the interface of NLP and computer vision. End-to-end SLMT between SL video and text in a spoken language exists, but performs poorly compared to translation broken down into intermediate steps where signs are represented by some orthographic form (e.g. in glosses (De Coster et al., 2023), or SL notation system (Walsh et al., 2022; Jiang et al., 2023)) - except for restricted domains and datasets (e.g. (Camgöz et al., 2020; Albanie et al., 2021; Zhang et al., 2023)). These subtasks are neatly shown in a diagram from Moryossef and Goldberg’s (2021) overview of the field in Figure 1.

Even though both text-to-SL gloss (Text2Gloss) and SL gloss-to-text (Gloss2Text) are sequence-to-sequence tasks using machine-readable text, the amount of parallel data available for any SL is orders of magnitude smaller than equivalent pairs of spoken languages. According to Duarte and colleague’s survey (2021), the largest parallel corpus between SL glosses and text available to researchers¹ contains 21,000 parallel utterances (Zhou et al., 2021a). It is reasonable to refer to all SLs as *extremely* low resource languages (Moryossef et al., 2021), and therefore data augmentation approaches must be adopted in order to improve the performance of translation models which include them.

In this paper, a novel method to generate semantic representations for Sign Language (SL) glosses is described. They are created by bootstrapping pre-trained word embedding models from spoken languages which already exist and their use is demonstrated in multilingual Text2Gloss and Gloss2Text machine translation experimental settings.

This paper is structured as follows: In Section 2, previous work where linguistic information is used to supplement gloss representations and its impact on SLMT is described, as well as work to create computational semantic resources for SLs in general. Section 3 sets out the process to generate SL gloss embeddings from pre-trained word embeddings, before Section 4 demonstrates their use in translation experiments. Findings from these experiments are described in Section 5 and discussed in Sections 6 and 7, along with potential future research directions using these embedding representations.

2 Background

One way of mitigating the semantic bottleneck created by gloss representation of signs is to explore techniques for low-resource neural machine translation (Sennrich and Haddow, 2016). These include data augmentation methods involving linguistic features (Armengol Estapé and Ruiz Costa-Jussà, 2021) as well as techniques specifically designed for Text2Gloss translation (Moryossef et al., 2021; Zhou et al., 2021b).

Zhu and colleagues’ (2023) comprehensive

study of these methods found, for DGS² corpora, that: (1) a combination of data augmentation strategies, and (2) transfer learning³ are viable methods to improve translation performance for Text2Gloss. They also highlight that it is important to ensure that these findings are generalisable to other SLs so further investigation such as the present study is required.

Other studies focused on injecting linguistic features into the embedding table for Text2Gloss and Gloss2Text (Egea Gómez et al., 2022; Chiruzzo et al., 2022), and found that transfer learning was again beneficial, as well as using linguistic features such as part-of-speech (PoS) and syntactic dependency tags - including PoS tags for SL glosses (McGill et al., 2023).

It may also be beneficial to use semantic information about signs into translation models, instead of (or as well as) using syntactic or grammatical information. No previous study with a parallel methodology was found, but other studies do use embeddings as part of SLMT models. Walsh and colleagues (2022) use sentence-level word2vec (Mikolov et al., 2013) or BERT (Devlin et al., 2019) embeddings to support Text2Notation (in HamNoSys (Hanke, 2004)) translation. Other studies use visual embeddings to support joint Sign2Gloss2Text (De Coster et al., 2023), SL recognition (Wong et al., 2023), or to encode phonological information for isolated SL recognition (Kezar et al., 2023).

This paper investigates using semantic information about words and glosses as a transfer learning strategy, and also its performance in combination with the syntax-based data augmentation methods seen in previous works.

2.1 Semantic data sources

Despite the fact that there are many word embedding collections for a great number of spoken languages, the same cannot be said about SLs. This problem is accentuated because the size of current SL corpora is not large enough to create high quality word embedding sets. Schuurman and colleagues (2023) propose SignNets, a database containing rich information about signs in a given SL, indexed by either gloss or an equivalent lexeme in a spoken language. This type of representation would be ideal to map meaning between signs and

¹How2Sign intended to include 35,000 parallel English/ASL text/glosses, but annotation was suspended indefinitely.

²German Sign Language (Deutsche Gebärdensprache)

³Such as pre-training on larger, language-agnostic models

between SLs and spoken language senses. However this research is in its early stages, and not ready for use in applications such as SLMT yet.

In contrast, Signbanks (Cassidy et al., 2018) are a well-established and extensible lexicon resource. Signbanks typically store information like ID-glosses, definitions or equivalent senses in a spoken language, phonological specification, images, and video for a given sign.

Semantic resources which allow the understanding of meaning in context, or calculating similarity of a given lexeme to another, are known as pre-trained word embedding models. There are no extant models of this type for SLs, which means that novel ones must be created. However, training models like word2vec or GloVe (Pennington et al., 2014) requires a large quantity of written utterances and it has been established that written SL data does not exist anywhere in large quantities.

Fortunately, it is possible to leverage data from the ambient spoken language in which glosses are written: For example, English for Auslan glosses or Dutch for glosses in Flemish Sign Language. One possible approach could be to just use pre-trained word embeddings without any modification for SL data - *e.g.* a Spanish word2vec model for Spanish Sign Language (LSE)⁴ tasks. However, in previous studies this approach has been shown to degrade the performance of Gloss2Text (Chiruzzo et al., 2022) and PoS-tagging (McGill et al., 2023) tasks.

Moreover, in studies in spoken languages, it has been shown that using high-quality English pre-trained embeddings as “anchors” to train bilingual word embedding models for low-resource languages is a promising strategy (Eder et al., 2021). Another study shows that English-lower resource languages bilingual lexica can be used to bootstrap the development of NLP-based tools in under-resourced languages (Wang et al., 2022).

3 Sign Language gloss embeddings

The motivation behind the present methodology is the proposition that, by mapping ID-glosses and their equivalent senses from a Signbank, it is possible to alter the weights of a pre-trained word embedding model from a spoken language in order to simulate the semantic interactions between signs in a given SL.

⁴Lengua de Signos Española

Spoken Language			Sign Language		
ID	Dims.	#Embs	ID	#Signs	#Embs
English	300	3.00M	ASL	5079	+2605
Finnish	100	247k	FinSL	3120	+1178
Spanish	300	1.00M	LSE	1221	+316
Dutch	320	627k	NGT	4144	+2938

Table 1: Left hand side: For each spoken language, the dimensionality and total number of word embeddings in its word2vec model. Right hand side: For each SL, number of signs in its Signbank(s) and the number of additional word vectors added to the new, bootstrapped word2vec model

As such, each gloss in a given SL is mapped to pre-trained embedding weights in one of three ways, along with some examples from the SLs in this study:

1. If the mapping between ID gloss and spoken language senses in a given Signbank is **one-to-one**, use the embedding weights from that sense (usually this is the same lexical item *e.g.* “TIME-D”⁵ = “time” in NGT⁶)
2. If there is a **one-to-many** relationship between ID gloss and spoken language senses, take the mean embedding weight from those senses (*e.g.* “WATER” ∈ {“water”, “to drink”} in LSE, “RAT” ∈ {“rat”, “rodent”, “mouse”, “freshman”, “rookie”} in ASL)
3. If there are no senses which match any existing sense in the spoken language Signbank, use the embedding weight from the word embedding model’s ‘unknown’ token (*e.g.* “GALLAUDET”⁷ = “UNK” in ASL)

The rest of the weights in the original pre-trained word embedding models remain the same if there is no gloss with the same label and are retained in the model to allow for the mapping of out-of-vocabulary token mapping.

3.1 Datasets

In order to create these bootstrapped pre-trained word embedding models for glosses, a given SL/spoken language pair must have all of the following dataset types available: (1) a Signbank with ID-glosses and translations in the ambient spoken language, (2) a pre-trained word embedding model for the spoken language, (3) parallel corpora of

⁵Glosses derived from other languages than English are translated here

⁶Nederlandse Gebarentaal

⁷The name of a well-known University for DHH students

continuous signing utterances, with both text in the spoken languages and glosses as annotations.

As seen in SL resource surveys (Duarte et al., 2021; Moryossef and Goldberg, 2021), SL-spoken language pairs reaching all these criteria are few in number. Therefore, the bootstrapping of word embedding models and translation experiments are performed on the following language pairs: Spanish Sign Language (LSE)-Spanish; American Sign Language (ASL)-English; Sign Language of the Netherlands (NGT)-Dutch; and Finnish Sign Language (FinSL)⁸-Finnish.

In all language pairs, a word2vec model was chosen, then all unique gloss-definition pairs from a given Signbank were processed following the technique outlined at the beginning of Section 3. The Signbanks and pre-trained word embedding models chosen for each language pair are shown in Appendix B. Table 1 shows the resources used to generate gloss embeddings along with some statistics. The parallel corpora used for translation experiments and a description of data preprocessing is described in Section 4.

3.2 Examples

This section demonstrates the operation of the embedding creation methodology, and shows the potential effectiveness of more accurately representing semantic relations between SL glosses. What follows are examples of gloss embeddings, and then the embedding space is shown visually.

Using cosine similarity to obtain the most similar word vectors, it is possible to compare representations of the same gloss/word in NGT and Dutch respectively. For example, the meaning mapping for “STAGE⁹” ∈ {“theatre”, “the stage”, “stage acting”} results in a slightly different semantic field for NGT and Dutch respectively. In Dutch, the most similar words include “*theatre, folk theatre, play, Bolshoi*”. In NGT, the most similar (cosine similarity) glosses include “ACT-A, ACT-B, VIOLIN, play”, incorporating the verb senses of the gloss. Note that similar NGT glosses contain lexemes which only exist in NGT like “ACT-B”. This is a positive sign, as it shows it is possible to map semantic relations to novel lexemes.

Figure 2 shows a visual example, for “Africa” in English, and “AFRICA ∈ {Africa, continent, ge-



Figure 2: Top N similar words plot for “Africa” and “rat” in ASL (top) and English (bottom)

ography}” in ASL, as a 2D representation of the vector space (t-SNE (van der Maaten and Hinton, 2008)) and the twelve most similar lexemes for each, as well as for the “RAT” example from Section 3. In English, similar words tend to be the names of nations, whereas similar terms for ASL are more terms related to geographical features. For “rat”, the dominant sense seems to be related to the “rookie” definition in ASL, as opposed to the animal in English.

As seen in Table 1, some glosses introduced to the bootstrapped embedding models do not exist in the original spoken language models. An interesting example of this are three LSE glosses derived from the Spanish lexeme “blood”: “SANGRE1” ∈ {“passion”, “to carry sth. in the blood”}; “SANGRE2” = “blood”; “SANGRE3” ∈ {“glass”, “blood”}. A plot for this example is

⁸Suomalainen viitomakieli

⁹TONEEL, in NGT gloss

shown in Appendix C.

3.2.1 Vector space

Turning to an overview of the semantic space overall, Figure 3 was created by plotting these 300-dimensional vectors in joint 2D space (also with t-SNE): (1) all unique glosses from the LSE Signbank, which (2) have an entry in both the original Spanish, and bootstrapped LSE word2vec models. This plot shows that the vector space is altered by the transformations made by the present methodology, and hopefully means that the bootstrapped word2vec model can better simulate SL semantics.

However, it is important to note that the total of 1221 LSE gloss vectors plotted here are the only ones whose weights may have been altered, while the rest of the 1M vectors in the original Spanish word2vec remain the same. This is hopefully not a large concern, as one would expect glosses in a parallel corpus to largely overlap with the ones used to create the modified word2vec model (see also, Table 3 for statistics on overlap).

4 Translation experiments

For each language pair, we perform both Text2Gloss and Gloss2Text experiments in two settings. Firstly, a baseline (Section 4.1) with each parallel corpus for each language. Then, following previous similar experimental setups (Moryossef et al., 2021; Chiruzzo et al., 2022; Zhu et al., 2023), a *warm start* transfer learning approach (Section 4.2) is executed. In other words, first a translation model is pretrained with a larger silver corpus and shared silver and gold vocabulary, and then finetuned on the same parallel data as the baseline.

Translation experiments are performed using OpenNMT-py 3.4.2 (Klein et al., 2017). OpenNMT is an open source translation toolkit which is based on LSTM encoder-decoder model with attention. All other running parameters are set to default, unless stated in Appendix A.

4.1 Baseline experiments

The baseline experiments involve Text2Gloss and Gloss2Text translation between the four spoken language-SL pairs. The specific parallel (or ‘gold’) datasets are described in Section 4.3. In order to evaluate the utility of these novel word embedding representations in real translation settings, the encoder and decoder (or both) embedding spaces,

that start in a random state by default in OpenNMT, are replaced by our collections of word2vec embeddings. For example, in the Gloss2Text setting for NGT→Dutch, there are four experimental settings:

1. *Baseline* (= default OpenNMT encoder/decoder parameters)
2. *Baseline-enc* (= NGT word2vec model encoder, OpenNMT default decoder)
3. *Baseline-dec* (= OpenNMT default encoder, Dutch word2vec model decoder)
4. *Baseline-both* (= NGT word2vec model encoder, Dutch word2vec model decoder)

This repeated for each language, and in the Text2Gloss direction, results in a total of 32 baseline experiments. Each setting is repeated for three runs of 10k epochs, starting at a random seed.

4.2 Pretrain + finetune experiments

Like in Section 4.1, there are also four experimental types: *PT+FT*, *PT+FT-enc*, *PT+FT-dec*, and *PT+FT-both*. However, for these experiments, models are trained on a larger parallel ‘silver’ dataset which is comprised of utterances in a spoken language alongside *pseudo*-glosses created by rule-based methods of data augmentation (*c.f.* (Moryossef et al., 2021; Chiruzzo et al., 2022; Zhu et al., 2023)).

During the pretraining phase, models are trained for three runs of 10k epochs on the parallel silver data. This phase also follows a warm start strategy by means of joint vocabulary (Nguyen and Chiang, 2017), whereby vocabulary is generated at the start of pretraining containing all tokens from both the *silver* and *gold* datasets. From each run, the best-performing model (BLEU measured for models at every 200 steps, based on the dev set) is chosen. These models are then fine-tuned for a further 5k epochs (three runs each) on the parallel spoken language/S� corpora from the Baseline experiments.

4.3 Datasets

Owing to the way SL datasets are collected, along with their low resource nature, we adopt different strategies for: (a) creating the silver datasets for each language pair, and (b) doing dataset splits in the gold corpora.

This section describes, by language pair, the gold and silver parallel datasets used in translation

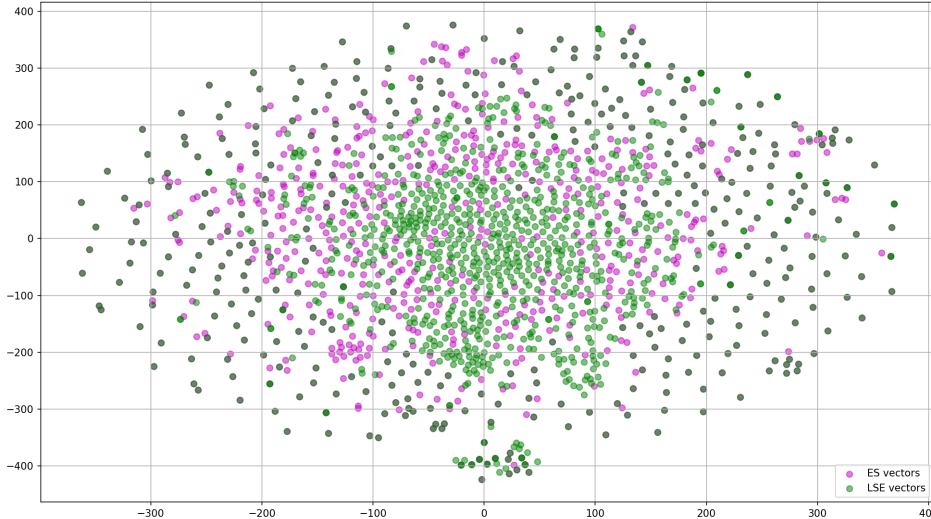


Figure 3: Spanish (purple) and LSE (green) word vectors from the LSE Signbank vocabulary plotted in joint 2d space. Grey points are where the weights are equal for both languages

	ASL/en	FinSL/fi	LSE/es	NGT/nl
Silver	87.7k	24.0k	20.3k	161k
Gold-train	2328	3480	1900	11.9k
Gold-dev	251	449	475	1484
Gold-test	352	534	482	1484
Gold-all	2931	4463	2857	14.8k

Table 2: Number of parallel utterances per language pair divided into dataset splits

	Baseline		PT+FT	
	#toks	overlap	#toks	overlap
ASL voc.	2410	75.8%	14.1k	73.7%
FinSL voc.	814	95.2%	2684	66.4%
LSE voc.	1123	60.7%	10.4k	83.5%
NGT voc.	3277	85.7%	25.2k	73.0%
en vocab	2377	95.1%	17.8k	78.6%
fi vocab	4523	24.4%	7450	24.2%
es vocab	2705	65.6%	16.4k	95.7%
nl vocab	11.2k	35.6%	38.0k	49.2%

Table 3: Vocabulary statistics for each language: number of unique tokens, and % overlap of tokens between the word2vec model and vocabulary in the gold (left columns) and silver+gold (right columns) datasets

experiments, dataset splits, and the methods used to generate silver data. Tables 2 and 3 show statistics about these datasets.

ASL/English: The NCSLGR and ASLLRP Corpora (Neidle et al., 2022) are combined as both datasets are relatively small for the present task. This data was accessed through ASLLRP’s Data Access Interface. These multimodal datasets contain utterances from twelve unique signers and contain a mixture of storytelling and elicited utterances, similar to the other parallel corpora used in this study. Like in Moryossef et al. (2021), the

silver data is the sample set¹⁰ from the ASLG-PC12 dataset - a parallel corpus where the ASL pseudo-glosses are generated with a linguistically-motivated rule-based approach. NCSLGR has been used before on its own in comparable studies (Zhu et al., 2023), but the decision was made to combine the two publicly-available glossed corpora so that as much parallel gold data as possible was available. The gold corpus was split into training-dev-test sets as close to 80%-10%-10% as possible, while also ensuring that the each unique signer only appears in one of these splits.

FinSL/Finnish: Corpus FinSL (Salonen et al., 2020) is used as the gold standard parallel dataset. For the silver data, Moryossef’s (Moryossef et al., 2021) language-agnostic rules for synthetic SL gloss generation is performed on 24k monolingual Finnish sentences selected at random (minimum 3 words per original utterances, duplicates removed) from the Tatoeba¹¹ collection. In addition, all first person pronouns are replaced with the gloss “OS:” (*pointing at self*) and other pronouns with “OS:minä” (*pointing sign*) to mirror the contents of the Corpus FinSL. This dataset was split 78%-10%-12% for train/dev/test.

LSE/Spanish: The iSignos Corpus from CORLSE (Cabeza and García-Miguel, 2019) is used for this language pair. There are 10 unique signers in this corpus, which informed the 64%-17%-19% train/dev/test split which is also used in previous studies (McGill et al., 2023). The silver

¹⁰https://huggingface.co/datasets/aslg_pc12

¹¹<https://tatoeba.org/en/>

data is also created using the same methodology from these studies, but using Tatoeba monolingual Spanish data to generate pseudo-glosses, and with slight differences in preprocessing decisions as described in Section 4.4.

NGT/Dutch: This language pair uses the largest parallel corpus available in this study, the CorpusNGT (Crasborn and Zwitterlood, 2008). Following SLMT experiments in the SignON project (Saggion et al., 2021), the dataset is split into partitions of 80%-10%-10%. Silver data was taken from a subset of the SONAR dataset for Dutch, and then modified with a rule-based approach (Bram Vanroy, *p.c.*) including gloss re-ordering¹² originally devised for Flemish Sign Language (VGT¹³).

4.4 Preprocessing

All four parallel corpora are annotated separately by dominant and non-dominant hand for SL glosses. As ML-based models, including NMT models, typically take linear alphanumeric input - it is necessary to modify the gloss annotations from these datasets. A systematic approach following *e.g.* Östling and colleagues (2017) was taken to linearise and lexicalise glosses:

- If two equal glosses occur simultaneously, only retain one
- If two different glosses simultaneously, place dominant hand gloss before non-dominant hand gloss
- Remove gestures which are not lexical signs
- Remove phonological features, tags indicating fingerspelling/name signs etc. from glosses

However, unlike similar studies which remove most affixes and labels, care was taken to match gloss labels in the parallel utterances to what is present in a given Signbank. In order to do this, glossing conventions and/or style guides such as SLAASh (Hochgesang, 2022) for ASLLRP and RADIS (Pérez et al., 2019) for CORLSE were referred to.

The same approach is taken for silver data generation. For example, pronouns which resemble

¹²<https://clin2022.uvt.nl/data-augmentation-for-machine-translation-of-sign-language-of-the-netherlands-and-flemish-sign-language/>

¹³Vlaamse Gebarental

those in the ambient spoken language, or where the silver dataset has its own gloss conventions, were edited to match what is used in the gold corpus/Signbank. In the synthetic ASL, all adjectives contained the prefix “DESC-”. As this does not occur in the gold data, they were removed. All gloss and spoken language text data is tokenised and in lowercase.

4.5 Evaluation

The best models from all runs of each experimental setting are evaluated on the held-out test set in the following way:

BLEU (Papineni et al., 2002) and **CHrF** (Popović, 2015) are the primary means of automatic evaluation in this study, measured using **sacreBLEU** (Post, 2018). BLEU-4 is calculated with disabled internal tokenisation¹⁴ (Müller et al., 2023). **METEOR** (Banerjee and Lavie, 2005) is also calculated through **nlTK**¹⁵ and reported. As there are three runs per experimental setup, **mean and standard deviation** are reported.

Statistical significance testing is also performed by means of **paired bootstrap resampling** (Koehn, 2004) calculated with Graham Neubig’s script¹⁶. Koehn states that this method of calculating significance at a level of $p < 0.05$ is effective with test sets greater than $N=300$. In this study, all test sets range between $N=352$ and $N=1484$.

Some **qualitative evaluation** is provided in the form of perceptive comments by the authors. Qualitative evaluation is of utmost importance to MT as a field¹⁷, especially low-resource MT where output with reasonable BLEU scores may still be ungrammatical or incomprehensible to the reader. Unfortunately, it was beyond the scope of this study to provide a more formal approach to qualitative assessment such as Direct Assessment (Graham et al., 2017; Zhu et al., 2023).

4.6 Reproducibility

The data, experimental configuration files, preprocessing and data augmentation scripts, scripts to generate embeddings, and model outputs for testing are all openly available¹⁸ for the purposes of transparency and reproducibility.

¹⁴Signature:

¹⁵<https://www.nltk.org/>

¹⁶<https://github.com/neubig/util-scripts/blob/master/paired-bootstrap.py>

¹⁷<https://bricksdont.github.io/posts/2020/12/seven-recommendations-for-mt-evaluation/>

¹⁸https://github.com/euan-mcgill/gloss_embeddings

5 Results and analysis

Table 4 summarises the quantitative findings of this study, reporting the best-performing model for each setup. Table 5 in Appendix D reports the best model on average (and standard dev.) across three runs for each setup. For the experimental setup acronyms used in this section and Table 4, refer to their descriptions in Sections 4.1 and 4.2.

For **es**→**LSE Text2Gloss**, the Baseline models with any kind of embeddings improved over the baseline in CHRf and METEOR, but only *PT+FT-both* performed better on the BLEU metric and this difference was not significant ($p = 0.25$, $N=482$). All *PT+FT* conditions had significantly higher BLEU scores than the Baseline. Within the *PT+FT* experimental setups, all metrics were markedly higher in the embedding setups, and *PT+FT-enc* ($p = 0.03$, $N=482$) and *PT+FT-both* ($p = 0.03$, $N=482$) showed a significant improvement. For **LSE**→**es Gloss2Text**, *PT+FT* tends to be a better strategy with *PT+FT-enc* being the only setup which performs significantly better than the baseline ($p = 0.03$, $N=482$), and higher scores in both metrics.

The fact that *PT+FT-both* performs significantly better than *PT+FT* in Text2Gloss, and *PT+FT-enc* than *Baseline* in Gloss2Text, is particularly promising as these conditions include the bootstrapped word embedding models for LSE.

For **nl**→**NGT Text2Gloss**, using embeddings improves BLEU scores in all setups, but only *PT+FT-dec* (with NGT bootstrapped glosses) in METEOR as well as being the only significant improvement on BLEU ($p < 0.01$, $N=1484$). The results are the mirror image in *PT+FT*: All setups except *PT+FT-dec* significantly improve over the baseline, and *PT+FT-enc* with only Dutch word2vec embeddings improves over *PT+FT* ($p = 0.05$, $N=1484$). However, Table 5 indicates a marked degree of variance compared to other language pairs and setups. This would be interesting to investigate further.

In **NGT**→**nl Gloss2Text** word2vec embeddings, as well as pretraining and finetuning, seems to damage the performance of this translation direction. Across both of these language pairs, compared to the BLEU scores the METEOR scores are also quite weak (compare LSE and FinSL results). In this language pair in particular, Table 3 shows that there is a large disparity in size between a much larger Dutch vocab than NGT. Moreover, the

Dutch word2vec model has a very low token coverage with both the gold and silver+gold vocab used in these experiments (both less than 50%). The consequence of this may be that it is difficult to create links between the lexical items in both languages.

For **en**→**ASL Text2Gloss**, the use of word2vec embeddings improves performance on most settings on both metrics. In the *PT+FT* setting, encoder English embeddings and both English and ASL embeddings improve significantly over the baseline (*PT+FT-both*: $p < 0.00$, $N=352$). For **ASL**→**en Gloss2Text**, over the baseline, significant improvements are seen when ASL embeddings are used in the encoder to support glosses: *PT+FT-enc* ($p = 0.02$, $N=352$), and *PT+FT-both*: ($p < 0.01$, $N=352$). It is therefore reasonable to infer that: (a) richer semantic representations for ASL, and (b) *warm-start* transfer learning on a larger silver dataset and joint vocabulary provides a real boost to translation performance and generalisibility.

It may be the case that there are no marked improvements within *PT+FT* as the en/ASL test partition is the smallest between all four language pairs (see Table 2). Also, unusually among parallel corpora, ASL’s vocabulary is actually larger than the English one where there is usually a large disparity in the other direction (see Table 3).

In the ASL/en language pair in general, the METEOR scores are much stronger compared to the others in this study. Perhaps a more even total of unique tokens in both ASL and English lexica contributes to this.

As for **fi**→**FinSL Text2Gloss** and **FinSL**→**fi Gloss2Text**, *PT+FT* with silver data provides an improvement in metrics across the board. Significant improvements are only seen in Gloss2Text: FinSL embeddings significantly improve in *Baseline-enc* over the baseline ($p = 0.02$, $N=534$), and in *PT+FT-enc* over *PT+FT* ($p = 0.04$, $N=534$), but curiously not over the baseline which has a very low BLEU score < 1 .

The results for FinSL/fi may be considered rather unusual on the whole. It is possible that the very low vocabulary size of the FinSL gold utterances ($N=814$), the disparity¹⁹ between this and the Finnish vocabulary size ($N=4523$), the replacement of all pronouns with just two lexemes (see

¹⁹BLEU has a brevity penalty, so the short sentences output for FSL should contribute to low scores

Best models - Text2Gloss	LSE			NGT			ASL			FinnSL		
	BLEU	CHrF	Met.	BLEU	CHrF	Met.	BLEU	CHrF	Met.	BLEU	CHrF	Met.
Baseline	7.34	0.198	0.089	18.66	0.269	0.115	15.46	0.372	0.286	5.54	0.174	0.109
Baseline+enc	7.32	0.212	0.102	19.31	0.261	0.112	17.56	0.398	0.312	6.41	0.195	0.125
Baseline+dec	7.12	0.202	0.093	19.92*	0.271	0.118	15.98	0.363	0.283	4.88	0.166	0.103
Baseline+both	7.64	0.217	0.106	19.66	0.266	0.114	16.30	0.381	0.303	6.27	0.193	0.128
PT+FT	9.94*	0.240	0.151	22.34*	0.310	0.142	18.46*	0.423	0.344	7.06	0.222	0.150
PT+FT+enc	17.83*†	0.341	0.197	22.67*†	0.306	0.144	20.26*	0.432	0.349	7.09	0.250	0.174
PT+FT+dec	16.48*	0.309	0.184	19.75	0.307	0.139	18.73	0.409	0.331	7.13	0.224	0.150
PT+FT+both	18.15*†	0.347	0.198	21.70*	0.311	0.140	19.67*	0.436	0.354	7.58	0.253	0.177
Best models - Gloss2Text	LSE			NGT			ASL			FinnSL		
	BLEU	CHrF	Met.	BLEU	CHrF	Met.	BLEU	CHrF	Met.	BLEU	CHrF	Met.
Baseline	7.80	0.193	0.146	4.84	0.220	0.144	14.29	0.352	0.356	0.90	0.116	0.119
Baseline+enc	8.25	0.192	0.159	4.47	0.219	0.144	13.61	0.353	0.357	1.27*	0.128	0.122
Baseline+dec	7.09	0.181	0.147	4.42	0.217	0.143	13.78	0.347	0.345	1.27	0.132	0.116
Baseline+both	8.70	0.197	0.161	4.31	0.214	0.140	13.25	0.347	0.352	1.75	0.151	0.131
PT+FT	8.67	0.201	0.168	3.38	0.200	0.129	16.52*	0.410	0.422	2.30*	0.164	0.145
PT+FT+enc	9.64*	0.211	0.178	3.59	0.208	0.132	16.88*	0.419	0.425	3.14†	0.177	0.155
PT+FT+dec	7.95	0.212	0.165	3.55	0.202	0.128	15.94	0.398	0.412	2.48†	0.164	0.138
PT+FT+both	9.02	0.214	0.179	3.53	0.206	0.131	17.05*	0.421	0.424	3.05	0.177	0.150

Table 4: Results summary for translation experiments in OpenNMT (* = significantly better than Baseline, † = significantly better than PT+FT). For each metric (Met. = METEOR), a higher score implies better performance.

Section 4.3), and the low token coverage of the Finnish word2vec model of Finnish tokens (24% in both the gold and silver+gold datasets) contribute to these results. Besides this, the FinSL dataset contains a few signs corresponding to descriptive markers (Salonen et al., 2019) (e.g. “_kvkk” for ‘whole object’ and “_kvmk” for ‘shape and size’) that are frequent (around 19% of the total signs in the training set). These signs are not lexical and have no corresponding ambient language lexemes, so an “unknown” random embedding was assigned to them.

These tokens’ high frequency may also explain the low performance of the FinSL \leftrightarrow fi experiments. In the future, we want to explore the possibility of creating embeddings for these markers using the average embedding of their corresponding Finnish descriptions.

5.1 Qualitative analysis

After a high-level comparison of model output which uses word vectors from this study, it is possible to observe lexical differences across experimental settings. These include, especially in lower-performing language pairs like es \rightarrow LSE, the replacement of more similar glosses even when the translation is inaccurate, a more similar distribution of PoS categories compared to the gold translation, and a lower prevalence of garbled output and model hallucination in lower-performing language pairs. Some qualitative examples from experimental settings are shown in Appendix E.

Looking at model output utterances, in tandem

BLEU Score	Interpretation
< 10	Almost useless
10 - 19	Hard to get the gist
20 - 29	The gist is clear, but has significant grammatical errors
30 - 40	Understandable to good translations
40 - 50	High quality translations
50 - 60	Very high quality, adequate, and fluent translations
> 60	Quality often better than human

Figure 4: Interpretability of BLEU scores

with the low BLEU scores, may explain unusual patterns of significance for FinSL/fi experiments. Figure 4 is a BLEU interpretability chart²⁰ which is useful to refer to when interpreting the quantitative results.

6 Discussion and limitations

So far, this exploratory work shows that using semantic representations tailored to SLs (in this case word2vec embeddings adapted to particular SL settings) is a promising avenue of research. Overall, the results present a positive outlook concerning the effectiveness of including bootstrapped word embedding models in the encoder and/or decoder of OpenNMT for Text2Gloss and Gloss2Text translation. In all *PT+FT*-* settings, the use of embeddings improved translation performance in at least one setting. This is also true in the baseline setting, apart from with NGT \rightarrow nl and

²⁰<https://cloud.google.com/translate/docs/advanced/automl-evaluate>

ASL→en. Many of these improvements were significant, and those where *PT+FT-embedding* significantly improved against *PT+FT* are particularly notable.

However, it is necessary to examine more data augmentation methods, types/sizes of word embedding models, sub-word tokenisation, and techniques to adapt semantic representations for the *extremely* low-resource setting of SL processing. It may also be worthwhile to attempt this approach on low-resource pairs of spoken languages, especially those with little or no written data (Aeppli et al., 2023) as anchor word embeddings already exist for spoken languages (Eder et al., 2021). Other practical tasks involving word embedding model support may include the tagging and parsing of SL gloss data (Östling et al., 2017; Yang and Zhang, 2018; García-Miguel and Cabeza, 2020).

Besides the use of OpenNMT for experiments, trying alternative open source translation toolkits such as MarianMT (Junczys-Dowmunt et al., 2018) (such as Perea-Trigo et al. (2024) for LSE).

Pretrained models like mBART (Liu et al., 2020) could also be a fruitful direction of research. Some preliminary experiments following Egea Gómez and colleagues (2022) were also attempted, using a mBART translation approach for LSE↔Spanish. However, some issues were found when applying the present method to mBART: Firstly, the model uses SentencePiece tokenisation, while this study’s embeddings are created with simple whitespace tokenisation. Furthermore, the mBART model expects a unified embedding space between source and target languages, which could skew the results for glosses that have the same surface form as ambient language words. It is possible to overcome these limitations, but given time and resource constraints the mBART experiments remain out of the scope of this work, and it is planned to explore them further in the future.

It would also be rewarding to explore other lexical SL resources such as Signpuddle²¹ which has been used in work on Text2Notation (Jiang et al., 2023) translation work. In addition, when SignNets (Schuurman et al., 2023) are further developed and contain rich metalinguistic information for many SLs, these will be a crucial resource for further studies in this area.

Some researchers may disagree with the use of glosses as a representation in SL processing

²¹<https://www.signbank.org/signpuddle/>

altogether, and disprefer splitting SLMT into a pipeline of intermediate tasks instead of treating it as an end-to-end task (Yin and Read, 2020). This is a valid position, and other work involving semantic representations in, for example, Video2Text could be complementary to studies like the present one.

Recent innovations into data-intensive methods such as 0-shot MT and NLP tasks often exclude SLs, because even though messy, unorganised, and seemingly irrelevant text data can be used for tasks in many spoken languages, this is not necessarily the case for the multimodal nature of SLs (Yin et al., 2021; Núñez-Marcos et al., 2023). However, recent research into *true 0-shot* translation; using LLMs to read and interpret reference material about the grammar of a language (Tanzer et al., 2024) - may aid SLMT and SL processing beyond that.

The large amount of experimental settings and limited computing resources available also meant that it was not possible to complete all of the evaluation that was initially planned. For example, from the insights gleaned from NGT→Dutch, it would be interesting to quantitatively investigate the connection between word embedding model’s vocab coverage and model performance. Qualitative analysis, though present, was unfortunately minimal and not formal and it would be greatly beneficial to expand it.

7 Concluding remarks

This study cast a wide net in order to devise novel methods to create semantic representations for SL glosses, and test their effectiveness when being used in SLMT. These experiments showed mixed but overall positive results, whereby bootstrapped pre-trained word embeddings from a spoken language can be modified with the present methodology in order to represent the semantic relations between SL glosses. It also provides further evidence that pretraining on silver data is effective across language pairs.

Future work will benefit from further experimentation with the methods undertaken to generate vector representation for signs whether represented by gloss, SL notation system, pose, or video frame. These embedding representations sit at the interface of NLP and computer vision-based approaches to SLMT, and characterise the need to follow both avenues of this field of research in a complimentary manner.

Acknowledgements

Amb el suport del Departament de Recerca i Universitats de la Generalitat de Catalunya.

This work is part of Maria de Maeztu Units of Excellence Programme CEX2021-001195-M, funded by MICIU/AEI /10.13039/501100011033

This work was supported in part by Oracle Cloud credits and related resources provided by Oracle for Research.

This work is a continuation of research started within the SignON project. SignON is a Horizon 2020 project, funded under the Horizon 2020 program ICT-57-2020 - “An empowering, inclusive, Next Generation Internet” with Grant Agreement number 101017255.

References

- Aepli, Noëmi, Chantal Amrhein, Florian Schottmann, and Rico Sennrich. 2023. A benchmark for evaluating machine translation metrics on dialects without standard orthography.
- Albanie, Samuel, Gül Varol, Liliane Momeni, Hannah Bull, Triantafyllos Afouras, Himel Chowdhury, Neil Fox, Bencie Woll, Rob Cooper, Andrew McParland, and Andrew Zisserman. 2021. BOBSL: BBC-Oxford British Sign Language Dataset.
- Armengol Estapé, Jordi and Marta Ruiz Costa-Jussà. 2021. Semantic and syntactic information for neural machine translation: Injecting features to the transformer. *Machine Translation*, 35:3:3–17.
- Banerjee, Satantjeev and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In Goldstein, Jade, Alon Lavie, Chin-Yew Lin, and Clare Voss, editors, *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan, June. Association for Computational Linguistics.
- Bragg, Danielle, Oscar Koller, Mary Bellard, Larwan Berke, Patrick Boudreault, Annelies Braffort, Naomi Caselli, Matt Huenerfauth, Hernisa Kacorri, Tessa Verhoef, Christian Vogler, and Meredith Ringel Morris. 2019. Sign Language Recognition, Generation, and Translation: An Interdisciplinary Perspective. In *The 21st International ACM SIGACCESS Conference on Computers and Accessibility, ASSETS '19*, page 16–31, New York, NY, USA. Association for Computing Machinery.
- Cabeza, Carmen and José M. García-Miguel. 2019. iSignos: Interfaz de datos de Lengua de Signos Española (versión 1.0).
- Camgöz, Necati Cihan, Oscar Koller, Simon Hadfield, and Richard Bowden. 2020. Sign language transformers: Joint end-to-end sign language recognition and translation. In *CVPR 2020*, pages 10020–10030.
- Cassidy, Steve, Onno Crasborn, Henri Nieminen, Wessel Stoop, Micha Hulsbosch, Susan Even, Erwin Komen, and Trevor Johnston. 2018. Signbank: Software to support web based dictionaries of sign language. In Calzolari, Nicoletta, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Koiti Hasida, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asuncion Moreno, Jan Odijk, Stelios Piperidis, and Takenobu Tokunaga, editors, *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May. European Language Resources Association (ELRA).
- Chiruzzo, Luis, Euan McGill, Santiago Egea-Gómez, and Horacio Saggion. 2022. Translating Spanish into Spanish Sign Language: Combining rules and data-driven approaches. In Ojha, Atul Kr., Chao-Hong Liu, Ekaterina Vylomova, Jade Abbott, Jonathan Washington, Nathaniel Oco, Tommi A Pirinen, Valentin Malykh, Varvara Logacheva, and Xiaobing Zhao, editors, *Proceedings of the Fifth Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2022)*, pages 75–83, Gyeongju, Republic of Korea, October. Association for Computational Linguistics.
- Cormier, Kearsy, Onno Crasborn, and Richard Bank. 2016. Digging into signs: Emerging annotation standards for sign language corpora. In Efthimiou, Eleni, Stavroula-Evita Fotinea, Thomas Hanke, Julie A. Hochgesang, Jette Kristoffersen, and Johanna Mesch, editors, *Proceedings of the LREC2016 7th Workshop on the Representation and Processing of Sign Languages: Corpus Mining*, pages 35–40, Portorož, Slovenia, May. European Language Resources Association (ELRA).
- Crasborn, Onno A and IEP Zwitserlood. 2008. The corpus ngt: an online corpus for professionals and laymen.
- Crasborn, Onno, Richard Bank, Inge Zwitserlood, Els van der Kooij, Ellen Ormel, Johan Ros, Anique Schüller, Anne de Meijer, Merel van Zuilen, Yasmine Ellen Nauta, Frouke van Winsum, and Max Vonk. 2020. NGT dataset in Global Signbank.
- De Coster, Mathieu, Dimitar Shterionov, Mieke Van Herreweghe, and Joni Dambre. 2023. Machine translation from signed to spoken languages: State of the art and challenges. *Universal Access in the Information Society*, pages 1–27.
- De Sisto, Mirella, Vincent Vandeghinste, Santiago Egea Gómez, Mathieu De Coster, Dimitar Shterionov, and Horacio Saggion. 2022. Challenges with sign language datasets for sign language recognition and translation. In *Proceedings of the Thirteenth Language Resources and Evaluation Confer-*

- ence, pages 2478–2487, Marseille, France, June. European Language Resources Association.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Burstein, Jill, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Duarte, Amanda, Shruti Palaskar, Lucas Ventura, Deepti Ghadiyaram, Kenneth DeHaan, Florian Metz, Jordi Torres, and Xavier Giro-i Nieto. 2021. How2sign: A large-scale multimodal dataset for continuous american sign language. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Duquenne, Paul-Ambroise, Holger Schwenk, and Benoît Sagot. 2023. Sonar: Sentence-level multimodal and language-agnostic representations.
- Eder, Tobias, Viktor Hangya, and Alexander Fraser. 2021. Anchor-based bilingual word embeddings for low-resource languages. In Zong, Chengqing, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 227–232, Online, August. Association for Computational Linguistics.
- Egea Gómez, Santiago, Luis Chiruzzo, Euan McGill, and Horacio Saggion. 2022. Linguistically enhanced text to sign gloss machine translation. In *Natural Language Processing and Information Systems: 27th International Conference on Applications of Natural Language to Information Systems, NLDB 2022, Valencia, Spain, June 15–17, 2022, Proceedings*, pages 172–183.
- García-Miguel, José M. and Carmen Cabeza. 2020. Hacia un treebank de dependencias para la lse. *Hesperia: Anuario de Filología Hispánica*, 22:111–143, mar.
- Graham, Yvette, Qingsong Ma, Timothy Baldwin, Qun Liu, Carla Parra, and Carolina Scarton. 2017. Improving evaluation of document-level machine translation quality estimation. In Lapata, Mirella, Phil Blunsom, and Alexander Koller, editors, *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 356–361, Valencia, Spain, April. Association for Computational Linguistics.
- Grushkin, Donald A. 2017. Writing signed languages: What for? what form? *American annals of the deaf*, 161(5):509–527.
- Hanke, Thomas. 2004. Hamnosys—representing sign language data in language resources and language processing contexts. In *LREC 2004, WS on RPSLs*, pages 1–6, Paris, France.
- Hochgesang, Julie A. 2022. Slaash id glossing principles, asl signbank and annotation conventions, version 3.2.
- Jiang, Zifan, Amit Moryossef, Mathias Müller, and Sarah Ebling. 2023. Machine translation between spoken languages and signed languages represented in SignWriting. In Vlachos, Andreas and Isabelle Augenstein, editors, *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1706–1724, Dubrovnik, Croatia, May. Association for Computational Linguistics.
- Junczys-Dowmunt, Marcin, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia, July. Association for Computational Linguistics.
- Kezar, Lee, Jesse Thomason, and Zed Sehyr. 2023. Improving sign recognition with phonology. In Vlachos, Andreas and Isabelle Augenstein, editors, *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2732–2737, Dubrovnik, Croatia, May. Association for Computational Linguistics.
- Klein, Guillaume, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander Rush. 2017. OpenNMT: Open-source toolkit for neural machine translation. In Bansal, Mohit and Heng Ji, editors, *Proceedings of ACL 2017, System Demonstrations*, pages 67–72, Vancouver, Canada, July. Association for Computational Linguistics.
- Koehn, Philipp. 2004. Statistical significance tests for machine translation evaluation. In Lin, Dekang and Dekai Wu, editors, *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain, July. Association for Computational Linguistics.
- Liu, Yinhan, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.
- McGill, Euan, Luis Chiruzzo, Santiago Egea Gómez, and Horacio Saggion. 2023. Part-of-speech tagging Spanish Sign Language data and its applications in sign language machine translation. In Ilinykh, Nikolai, Felix Morger, Dana Dannélls, Simon Dobnik, Beáta Megyesi, and Joakim Nivre, editors, *Proceedings of the Second Workshop on Resources and Representations for Under-Resourced*

- Languages and Domains (RESOURCEFUL-2023)*, pages 70–76, Tórshavn, the Faroe Islands, May. Association for Computational Linguistics.
- Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space.
- Moryossef, Amit and Yoav Goldberg. 2021. Sign Language Processing. <https://sign-language-processing.github.io/>.
- Moryossef, Amit, Kayo Yin, Graham Neubig, and Yoav Goldberg. 2021. Data augmentation for sign language gloss translation. In Shterionov, Dimitar, editor, *Proceedings of the 1st International Workshop on Automatic Translation for Signed and Spoken Languages (AT4SSL)*, pages 1–11, Virtual, August. Association for Machine Translation in the Americas.
- Müller, Mathias, Zifan Jiang, Amit Moryossef, Annette Rios, and Sarah Ebling. 2023. Considerations for meaningful sign language machine translation based on glosses. In Rogers, Anna, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 682–693, Toronto, Canada, July. Association for Computational Linguistics.
- Neidle, Carol, Augustine Opoku, and Dimitris N. Metaxas. 2022. ASL video corpora & sign bank: Resources available through the american sign language linguistic research project (ASLLRP). *CoRR*, abs/2201.07899.
- Nguyen, Toan Q. and David Chiang. 2017. Transfer learning across low-resource, related languages for neural machine translation. In Kondrak, Greg and Taro Watanabe, editors, *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 296–301, Taipei, Taiwan, November. Asian Federation of Natural Language Processing.
- Núñez-Marcos, Adrián, Olatz Perez de Viñaspre, and Gorka Labaka. 2023. A survey on sign language machine translation. *Expert Systems with Applications*, 213:118993.
- Östling, Robert, Carl Börstell, Moa Gärdenfors, and Mats Wirén. 2017. Universal dependencies for swedish sign language. In *Nordic Conference of Computational Linguistics*.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Isabelle, Pierre, Eugene Charniak, and Dekang Lin, editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Pennington, Jeffrey, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Perea-Trigo, Marina, Celia Botella-López, Miguel Ángel Martínez-del Amor, Juan Antonio Álvarez García, Luis Miguel Soria-Morillo, and Juan José Vegas-Olmos. 2024. Synthetic corpus generation for deep learning-based translation of spanish sign language. *Sensors*, 24(5).
- Popović, Maja. 2015. chrF: character n-gram F-score for automatic MT evaluation. In Bojar, Ondřej, Rajan Chatterjee, Christian Federmann, Barry Haddow, Chris Hokamp, Matthias Huck, Varvara Logacheva, and Pavel Pecina, editors, *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal, September. Association for Computational Linguistics.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *3rd Conf. on MT*, pages 186–191, Belgium, Brussels. ACL.
- Pérez, Ania, José M. García-Miguel, and Carmen Cabeza. 2019. Corpus annotation for studying grammatical expression of events: notes about the design of radis project. *Sensos-e*, 6(1):40–61, Sep.
- Saggion, Horacio, Dimitar Shterionov, Gorka Labaka, Tim Van de Cruys, Vincent Vandeghinste, and Josep Blat. 2021. Signon: Bridging the gap between sign and spoken languages. In *Alkorta J, Gonzalez-Dios I, Atutxa A, Gojenola K, Martínez-Cámara E, Rodrigo A, Martínez P, editors. Proceedings of the Annual Conference of the Spanish Association for Natural Language Processing: Projects and Demonstrations (SEPLN-PD 2021) co-located with the Conference of the Spanish Society for Natural Language Processing (SEPLN 2021); 2021 Sep 21-24; Málaga, Spain. Aachen: CEUR Workshop Proceedings; 2021. p. 21-5.* CEUR Workshop Proceedings.
- Salonen, Juhana, Tuija Wainio, Antti Kronqvist, and Jarkko Keränen. 2019. Suomen viittomakielten korpusprojektin (cfinsl) annotointiohjeet. In *Annotation Convention. Helmikuu: Department of Linguistics and Communication Sciences, Sign Language Center, University of Jyväskylä*, page 40.
- Salonen, Juhana, Antti Kronqvist, and Tommi Jantunen. 2020. The corpus of Finnish Sign Language. In Efthimiou, Eleni, Stavroula-Evita Fotinea, Thomas Hanke, Julie A. Hochgesang, Jette Kristoffersen, and Johanna Mesch, editors, *Proceedings of the LREC2020 9th Workshop on the Representation and Processing of Sign Languages: Sign Language Resources in the Service of the Language Community, Technological Challenges and Application Perspectives*, pages 197–202, Marseille, France, May. European Language Resources Association (ELRA).

- Schuurman, Ineke, Thierry Declerck, Caro Brosens, Margot Janssens, Vincent Vandeghinste, and Bram Vanroy. 2023. Are there just WordNets or also SignNets? In Rigau, German, Francis Bond, and Alexandre Rademaker, editors, *Proceedings of the 12th Global Wordnet Conference*, pages 172–178, University of the Basque Country, Donostia - San Sebastian, Basque Country, January. Global Wordnet Association.
- Sennrich, Rico and Barry Haddow. 2016. Linguistic input features improve neural machine translation. In *1st Conf. on MT*, pages 83–91, Berlin, Germany. ACL.
- Tanzer, Garrett, Mirac Suzgun, Eline Visser, Dan Jurafsky, and Luke Melas-Kyriazi. 2024. A benchmark for learning to translate a new language from one grammar book.
- van der Maaten, Laurens and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605.
- Walsh, Harry, Ben Saunders, and Richard Bowden. 2022. Changing the representation: Examining language representation for neural sign language production. In Efthimiou, Eleni, Stavroula-Evita Fotinea, Thomas Hanke, John C. McDonald, Dimitar Shterionov, and Rosalee Wolfe, editors, *Proceedings of the 7th International Workshop on Sign Language Translation and Avatar Technology: The Junction of the Visual and the Textual: Challenges and Perspectives*, pages 117–124, Marseille, France, June. European Language Resources Association.
- Wang, Xinyi, Sebastian Ruder, and Graham Neubig. 2022. Expanding pretrained models to thousands more languages via lexicon-based adaptation. In Muresan, Smaranda, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 863–877, Dublin, Ireland, May. Association for Computational Linguistics.
- Wong, R., N. Camgoz, and R. Bowden. 2023. Learnt contrastive concept embeddings for sign recognition. In *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, pages 1937–1946, Los Alamitos, CA, USA, oct. IEEE Computer Society.
- Yang, Jie and Yue Zhang. 2018. Ncrf++: An open-source neural sequence labeling toolkit. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*.
- Yin, Kayo and Jesse Read. 2020. Better sign language translation with STMC-transformer. In Scott, Donia, Nuria Bel, and Chengqing Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5975–5989, Barcelona, Spain (Online), December. International Committee on Computational Linguistics.
- Yin, Kayo, Amit Moryossef, Julie Hochgesang, Yoav Goldberg, and Malihe Alikhani. 2021. Including signed languages in natural language processing. In Zong, Chengqing, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7347–7360, Online, August. Association for Computational Linguistics.
- Zhang, Biao, Mathias Müller, and Rico Sennrich. 2023. SLTUNET: A simple unified model for sign language translation. In *The Eleventh International Conference on Learning Representations*.
- Zhou, H., W. Zhou, W. Qi, J. Pu, and H. Li. 2021a. Improving sign language translation with monolingual data by sign back-translation. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1316–1325, Los Alamitos, CA, USA, jun. IEEE Computer Society.
- Zhou, Hao, Wengang Zhou, Weizhen Qi, Junfu Pu, and Houqiang Li. 2021b. Improving sign language translation with monolingual data by sign back-translation.
- Zhu, Dele, Vera Czehmann, and Eleftherios Avramidis. 2023. Neural machine translation methods for translating text to sign language glosses. In Rogers, Anna, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12523–12541, Toronto, Canada, July. Association for Computational Linguistics.

A OpenNMT parameters

(1) To **build** vocabulary:

```
python build-vocab.py -n_sample 50000
```

(2a) To **train** translation models (pre-training):

```
python train.py -feat_merge "concat" -
bucket_size 144 -world_size 1 -gpu_ranks [0]
-save_checkpoint_steps 200 -train_steps 10000
-valid_steps 200 -log_file "specified.log"
```

(2b) To **train** translation models (fine-tuning):

```
python train.py -feat_merge "concat" -
bucket_size 144 -world_size 1 -gpu_ranks [0]
-save_checkpoint_steps 200 -train_steps +5000
-valid_steps 200 -train_from "specified-pt-
model.pt" -reset_optim keep_states -log_file
"specified.log"
```

(3) To **translate** test data for evaluation:

```
python translate.py -ban_unk_token
```

B Signbanks and word2vec models used

For **ASL-English**, a combination of the ASL Signbank²² ASLLRP Sign Bank (Neidle et al., 2022) are used and the GoogleNews word2vec (Skipgram) model²³.

For **FinSL-Finnish**, it is the Suomen Signbank²⁴ and the Finnish Text Collection word2vec model²⁵.

For **LSE-Spanish**, the CORLSE lexicon gathered from the iSignos Corpus’ web resource (Cabeza and García-Miguel, 2019) as well as the Spanish Billion Words model²⁶.

And for **NGT-Dutch**, the Global Signbank (NGT dataset) (Crasborn et al., 2020) and SONAR embeddings (Duquenne et al., 2023).

C Vector similarity plots

Figure 5 shows the ten most similar (cosine similarity) word in the LSE word2vec model for the three glosses based on the lexeme “BLOOD” in LSE mentioned in Section 3.2, represented in 2D vector space.

D Results: Mean and standard deviation

Table 5 shows the best-performing model (number of training epochs shown) on average from three runs in each experimental setup. The *PT+FT* experiments only show one set of experimental runs, as recall that from the pre-training phase, the best-performing epoch from each of the three runs is chosen to fine-tune for another 5000 epochs on *gold* data. Results for FinSL \leftrightarrow fi could not be shown, as only one run per setup was undertaken.

Similar to the findings based on the best model in each setup shown in Table 4, for most language pairs *PT+FT* performed more strongly than the Baseline. Using features tends to improve translation results on average, but the standard deviation figures show a high degree of variance between settings, particularly when translating from Dutch.

E Qualitative analysis examples

This Appendix shows four utterances from different translation directions and experimental setups which exemplify the use of bootstrapped SL embeddings in the encoder or decoder.

Figure 6 is an example from es \rightarrow LSE. The original *gold* output sentence from the test set of iSignos is challenging, particularly as it contains a classifier predicate²⁷ “RECIBIR-MONTÓN”. Comparing the *PT+FT-both* hypothesis to *PT+FT*, notice that “child” is rendered more accurately as “HOMBRE PEQUEÑO2” rather than “HOMBRE PERSONA” (a frequent bigram in this corpus). Also, the first person plural pronoun is correctly identified. Whether or not having tailored semantic representations available to the decoder/SL output brings about this improvement is up for debate, but the output is more faithful to the *gold* output nonetheless.

As a counterexample, Figure 7 compares the *gold* output for the given sentence with the *Baseline*, *PT+FT*, and *PT+FT-both* hypotheses. In this case, it appears that *PT+FT* output reflects the semantics of ASL in a better way. The signs “GIVE” and “GIFT” are exemplars of the phenomenon in ASL where signs can be used as nouns, verbs, or adjectives interchangeably, so using either in this instance would be grammatical. As for the model using both word2vec embedding representations, it chooses “GO-OUT” which - while still a verb - would not necessarily be the best choice.

Finally, Figure 8 shows a more challenging example - again from Spanish \rightarrow LSE where no model can provide a grammatical output. The outputs from *Baseline* and *PT+FT* appear like model hallucinations of frequently-occurring tokens from the training data. The same may be said about the *PT+FT-both* output. However, the connection between “padres” and “PADRE”/“MADRE” appears to be more robust and appears in its hypothesis. The *PT+FT-both* hypothesis is the only one to include a negative “NO” (“NADA2” appears in the *gold* output) which may imply that using SL-derived embeddings may also be more robust to part-of-speech class.

²²<https://aslsignbank.haskins.yale.edu/>

²³<https://code.google.com/archive/p/word2vec/>

²⁴<https://signbank.csc.fi/>

²⁵<http://urn.fi/urn:nbn:fi:lb-2022041405>

²⁶<https://crscardellino.ar/SBWCE/>

²⁷Signs which are more iconic, which may be unique to a given signer, and do not have a fixed meaning *e.g* in a SL dictionary. These are used to depict or describe actions, entities, and events among other things.

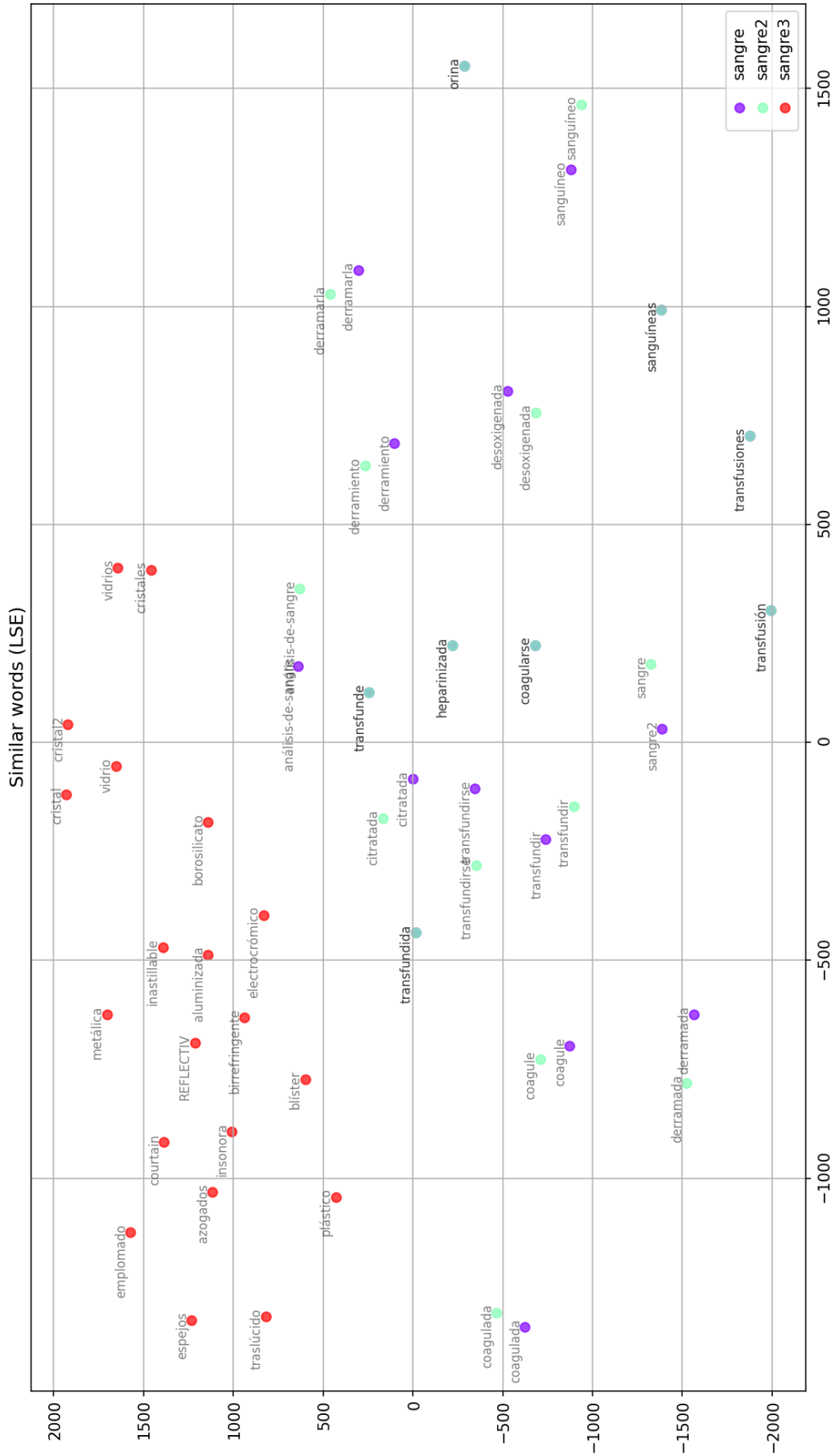


Figure 5: Top ten most similar lexemes to each gloss based on the Spanish word for “blood”

Mean + std. dev Text2Gloss	es→LSE		nl→NGT		en→ASL	
	Epoch	BLEU	Epoch	BLEU	Epoch	BLEU
Baseline	10000	5.55 ± 1.71	7600	12.63 ± 3.81	9000	14.74 ± 0.63
Baseline+enc	9600	5.47 ± 0.41	9400	11.54 ± 4.17	6400	16.94 ± 0.58
Baseline+dec	9600	4.30 ± 1.27	4600	16.92 ± 2.60	4800	14.30 ± 1.60
Baseline+both	10000	5.07 ± 3.15	8400	13.63 ± 0.90	7800	15.42 ± 0.33
PT+FT	7600+3200	9.12 ± 0.75	9200+1000	17.88 ± 2.34	1200+3200	18.01 ± 0.42
PT+FT+enc	7800+3000	16.50 ± 1.19	9800+3800	18.24 ± 2.11	6200+4400	18.84 ± 1.23
PT+FT+dec	5600+3400	15.40 ± 0.64	6200+1800	16.10 ± 4.64	6400+5000	17.92 ± 1.08
PT+FT+both	3800+4400	16.61 ± 0.47	6400+1600	18.20 ± 1.84	3200+4800	18.87 ± 0.69
Mean + std. dev Gloss2Text	LSE→es		NGT→nl		ASL→en	
	Epoch	BLEU	Epoch	BLEU	Epoch	BLEU
Baseline	3600	6.96 ± 0.73	7400	4.41 ± 0.43	5400	12.80 ± 1.42
Baseline+enc	3000	7.63 ± 0.29	7400	4.30 ± 0.15	4800	12.89 ± 0.41
Baseline+dec	4400	6.11 ± 0.59	9800	4.12 ± 0.18	5600	12.94 ± 0.75
Baseline+both	4000	7.75 ± 0.88	8400	4.10 ± 0.10	5400	12.68 ± 0.86
PT+FT	8600+2000	7.91 ± 0.69	4600+3200	3.18 ± 0.12	8800+1400	15.57 ± 0.54
PT+FT+enc	9800+3000	9.12 ± 0.53	9200+3200	3.37 ± 0.12	8400+3600	15.86 ± 0.59
PT+FT+dec	7800+1800	7.52 ± 0.29	5200+4000	3.09 ± 0.24	9400+1600	15.01 ± 0.76
PT+FT+both	8200+3400	8.49 ± 0.79	8400+4400	3.25 ± 0.11	9600+3800	16.53 ± 0.42

Table 5: Results summary for translation experiments in OpenNMT - BLEU-4 based mean and standard deviation for three runs in each experimental setup, along with the number of epochs for which the model is chosen. fi→FSL not shown as only underwent one run per setting.

Gold hyp: "HOMBRE PEQUEÑO2 REVISTA RECIBIR-MONTÓN INDX.PRO:1pl"
gloss: man small book receive-stack_{CL-M} we
ES: *Los niños nos dan los libros*

PT+FT hyp: "HOMBRE PERSONA REVISTA DAR INDX.PRO:3pl"
gloss: man person book give they

PT+FT-both hyp: "HOMBRE PEQUEÑO2 REVISTA DAR INDX.PRO:1pl"
gloss: man small book give we

Figure 6: Translation output from the Spanish sentence "The children give us the books" into LSE from the original corpus, and two model output hypotheses

Gold hyp: "PADRE^MADRE ENTENDER NADA2"
gloss: parents understand nothing
ES: *A mis padres no los entendía en absoluto*

Baseline hyp: "INDX.PRO:2sg PEQUEÑO2 CÓMO"
gloss: you small how

PT+FT hyp: "PROPIO HOMBRE ESPÍRITU"
gloss: own man mind

Gold hyp: "JOHN NOW i:GIVE:j CHOCOLATE MOTHERwg IX-3p:j"
gloss: John now _{REF1}give_{REF2} chocolate mother it
EN: *John is right now giving chocolate to mother.*

Baseline hyp: "JOHN NOW FINISH NOW SUE"
Gloss: John now already now Sue

PT+FT hyp: "JOHN NOW FINISH i:GIFT:j MOTHERwg"
gloss: John currently already _{REF1}give_{REF2} mother

PT+FT-both hyp: "JOHN NOW GO-OUT CHOCOLATE MOTHERwg"
gloss: John now go-out chocolate mother

Figure 7: Translation output from the English sentence "John is right now giving chocolate to mother" into ASL, and three model output hypotheses

PT+FT-both hyp: "INDX.PRO:1SG PADRE^MADRE NO"
gloss: I parents no

Figure 8: Translation model output from the Spanish sentence "As for my parents, I did not understand them at all", and three model hypotheses

Quality Estimation with k -nearest Neighbors and Automatic Evaluation for Model-specific Quality Estimation

Tu Anh Dinh¹

tu.dinh@kit.edu

Tobias Palzer²

tobiaspalzer.tp@gmail.com

Jan Niehues³

jan.niehues@kit.edu

^{1,3}Karlsruhe Institute of Technology
Karlsruhe, Germany

²Technical University of Munich
Munich, Germany

Abstract

Providing quality scores along with Machine Translation (MT) output, so-called reference-free Quality Estimation (QE), is crucial to inform users about the reliability of the translation. We propose a model-specific, unsupervised QE approach, termed k NN-QE, that extracts information from the MT model’s training data using k -nearest neighbors. Measuring the performance of model-specific QE is not straightforward, since they provide quality scores on their own MT output, thus cannot be evaluated using benchmark QE test sets containing human quality scores on premade MT output. Therefore, we propose an automatic evaluation method that uses quality scores from reference-based metrics as gold standard instead of human-generated ones. We are the first to conduct detailed analyses and conclude that this automatic method is sufficient, and the reference-based MetricX-23 is best for the task.

1 Introduction

Machine Translation (MT), due to its currently advanced stage in research, has been widely adopted in real-life use cases (Vieira et al., 2021). In many application domains such as health care or lawsuits, errors in translation could be tremendously harmful to the users. Therefore, it is important to

inform the user whether to rely on a certain translation, by providing some kind of quality assessment along with each translation output. This task is referred to as Quality Estimation (QE).

More specifically, Quality Estimation is assigning quality scores to MT output, without using gold-standard human translation. Common QE approaches train a standalone QE module that takes in the source sentences and the MT outputs to produce quality scores. These QE modules are usually model-agnostic, i.e., they can work with the output of any MT model. However, they often require training on human-labeled quality data, which can be costly to obtain. Another line of research is on model-specific QE, where they exploit or modify the MT model for self-quality assessment, thus not requiring training a separate QE module. Following this line of work, we propose k NN-QE - an unsupervised QE approach that exploits the information of inference-time output’s k -nearest neighbors found in the MT model’s training data. We hypothesize that the closer the inference-time sample output is to the training data, the better the quality of the translation, since it is an indication that the model has learnt about such samples. The QE scores obtained using our method can also be interpreted as the confidence scores of the MT model.

Unlike model-agnostic QE approaches which can take any MT translation as input, evaluating model-specific QE approaches like k NN-QE is not as straightforward. Public QE test sets are generated using human quality scores on pre-made MT output, thus not always suitable for QE approaches that perform self-evaluation on their own MT output by design. Many previous works on model-specific QE perform human evaluation on their own MT output to be used as gold standard to evaluate QE metrics (Riktors and Fishel, 2017;

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

^{1,2}Equal contributions.

²Tobias contributes during his Bachelor Thesis at KIT.

Niehues and Pham, 2019; Fomicheva et al., 2020; Zhang et al., 2022). However, for faster development, it would be useful to automatically evaluate QE metrics and not relying on human resource. Therefore, we propose using quality scores generated by reference-based metrics as the gold standard to automatically evaluate reference-free QE metrics. Our motivation is that reference-based metrics, by making use of reference translation, tend to be better than reference-free QE metrics (Freitag et al., 2023), thus can be used as gold standard. To the best of our knowledge, we are the first to perform a detailed analysis on whether reference-based metrics are sufficient to evaluate QE, and which reference-based metric is best suited. For this analysis, we make use of different QE submissions to public shared tasks: WMT22 Metrics (Freitag et al., 2022) and WMT22 Quality Estimation (Zerva et al., 2022). We investigate whether our automatic QE evaluation method can produce similar QE rankings compared to using human-labeled quality data in these shared tasks.

In summary, our contribution is in two folds:

1. A **model-specific, unsupervised QE** approach, termed $kNN-QE^*$, which exploits the similarity of MT generated output and MT models’ training data. Our main findings are: (1) $kNN-QE$ outperforms an unsupervised baseline using MT output probabilities, but falls behind supervised QE; and (2) $kNN-QE$ works with a small number of neighbors and partial access to the MT training data.
2. An **automatic QE evaluation** method[†] using a reference-based metric’s quality scores as gold-standard instead of human-labeled quality scores. Our main findings are: (1) QE ranking made by reference-based metrics correlate well with ones made by human quality scores; (2) Segment-level evaluation performance does not strictly correlate to QE ranking performance for reference-based metrics; and (3) MetricX 23 (Juraska et al., 2023) is the most robust for ranking QE metrics.

2 Related Work

Quality Estimation Quality Estimation (QE) aims to measure the quality of MT output without using human references. Common QE approaches

are model-agnostic, where a QE module takes in a source sentence, an MT translation and outputs a quality score (Blain et al., 2023). This approach has 2 drawbacks: (1) it requires a stand-alone module for QE, and (2) it requires human quality data to train the QE module, which can be costly.

Model-specific QE Researchers have also been looking into integrating Quality Estimation into MT models. These approaches exploit information or modify white-box MT model to measure the translation quality, rather than training a separate QE module relying completely on human quality data. (Rikters and Fishel, 2017) uses the attention distributions from the MT model as a QE metric. (Fomicheva et al., 2020) uses the attention distribution and the output probabilities from the MT model for QE. (Lu et al., 2022) propose QE learnt jointly with the training of the MT model. In their approach, the MT model can ask for hints to improve its translation, and the more hints it asks for, the lower the confidence. (Zhang et al., 2022) extends the MT model with a self-estimator module for QE, which examines whether it can reconstruct the source sentence’s semantics using the information from the decoding procedure. The work by (Niehues and Pham, 2019) is the closest to our $kNN-QE$, where they measure the similarity of the test sentence with sentences from the training data to estimate translation quality. The difference between this work and ours is that they focus on evaluating source side rather than target side; they use encoder output similarity rather than decoder output similarity; and they do not analyze different metrics derived from the nearest neighbors. Evaluating these model-specific QE approaches is not straightforward, as will be discussed below.

Automatic QE evaluation The standard way to evaluate QE is to use some benchmark test sets, containing human quality scores on the output of some MT models. The QE scores are then compared against the human scores on these pre-made translations. This works mostly for model-agnostic QE, since they can evaluate any MT output. However, for the model-specific QE approaches like $kNN-QE$, which provide quality scores on their own MT output, there are no longer readily available human quality scores for QE evaluation. Previous works on model-specific QE address this issue differently. Some works use MT glass-box features for QE without changing the

^{*}<https://github.com/TuAnh23/auto-meta-eval-qe>

[†]<https://github.com/TuAnh23/knn-box>

MT model, thus they can still produce the same MT translation that is used in the QE benchmarks (Yankovskaya et al., 2018; Wang et al., 2021). (Lu et al., 2022) train their MT model on the same data as the model used in the QE benchmarks and perform force decoding to get the exact same MT output. These approaches are then limited to the MT model used in the QE benchmarks. On the other hand, some works perform human evaluation on their own MT output for QE evaluation (Rikters and Fishel, 2017; Niehues and Pham, 2019; Fomicheva et al., 2020; Zhang et al., 2022). This requires human resource, which is costly and not always available. Overall, it is not yet clear what is the go-to method to perform automatic evaluation for model-specific QE. To the best of our knowledge, we are the first to perform detailed analysis on whether it is possible to automate evaluation for QE by making use of reference-based metrics.

***k*NN for generation tasks** Previous works have applied *k*-nearest neighbors in text generation. *k*NN-LMs (Khandelwal et al., 2019) enable language models to interpolate their token prediction output with a *k*-nearest neighbors model, where nearest neighbors are retrieved from a datastore of sample representations. *k*NN-MT (Khandelwal et al., 2020) also enables the MT model to predict tokens using a nearest neighbor classifier over a datastore of representations. *k*NN-LMs and *k*NN-MT are particularly useful for adapting models to diverse domains by using domain-specific datastores. Our *k*NN-QE approach is similar to these works in two aspects. First, in the datastore generation process, it also generates token representations by performing one forward pass of the model through the training data. Second, during inference, it also retrieves similar tokens in the datastore based on the token representation distance. The difference is that our *k*NN-QE approach uses the retrieved neighbors to assess the quality of the generated token, rather than modifying the model output like *k*NN-LMs and *k*NN-MT.

3 Quality Estimation with *k*NN

Motivation We propose *k*NN-QE - a model-specific Quality Estimation method that exploits information from the MT model’s training data using *k*-nearest neighbors. Our method is unsupervised, thus does not require human quality scores for training. Generally, if the hidden representation of a translation sample generated during in-

ference is similar to ones generated on the training data, then it is an indication that this sample is in-domain, thus more likely to have higher quality.

Generating the datastore We generate a datastore on the MT training data as follows. We first use the MT model to perform translation on its training set with force decoding on the reference. That is, we give the model human reference translation prefixes as input at every time step to generate the next translation token. We save the last-layer decoder hidden representation of every output token to the datastore. We do forced decoding on the reference for datastore generation since it provides an indication of confidence: if during inference, the self-generated prefix translation is high-quality, it would better match the forced decoding condition where prefixes are gold translation, thus making the representation of the inference-time generated token closer to the ones in the datastore.

Formally, let the m^{th} training source sentence be $X^m = (x_1^m, x_2^m, \dots, x_{|X^m|}^m)$ and the m^{th} training reference target sentence be $\hat{Y}^m = (\hat{y}_1^m, \hat{y}_2^m, \dots, \hat{y}_{|\hat{Y}^m|}^m)$, where the element tokens are subwords. The last-layer decoder hidden representation of the output token at time step i with forced decoding on the reference is:

$$\hat{d}_i^m = Dec(E^m, (\hat{y}_1^m, \hat{y}_2^m, \dots, \hat{y}_{i-1}^m)) \quad (1)$$

where $E^m = Enc(X^m)$, Dec and Enc are the decoder and encoder functions respectively. We save to the datastore the \hat{d}_i^m representation for each output token \hat{y}_i^m in the training data.

Retrieving neighbors during inference During inference, for each generated token, we use its last-layer decoder hidden representation and find the *k*-nearest neighbors from the datastore. The neighbor retrieval can be highly optimized using toolkits like Faiss (Johnson et al., 2019), thus does not cost too much inference speed.

Formally, let the output sentence be $Y = (y_1, y_2, \dots, y_{|Y|})$. The last-layer decoder hidden representation of the output token y_j at time step j is (no forced decoding at inference time):

$$d_j = Dec(E, y_1, y_2, \dots, y_{j-1}) \quad (2)$$

where $E = Enc(X)$. We find the set N_j of *k*-nearest neighbors of y_j by:

$$N_j = \underset{m \in \text{train}, i \in 1..|\hat{Y}^m|}{\operatorname{argmin}^k} (L_2(d_j, \hat{d}_i^m)) \quad (3)$$

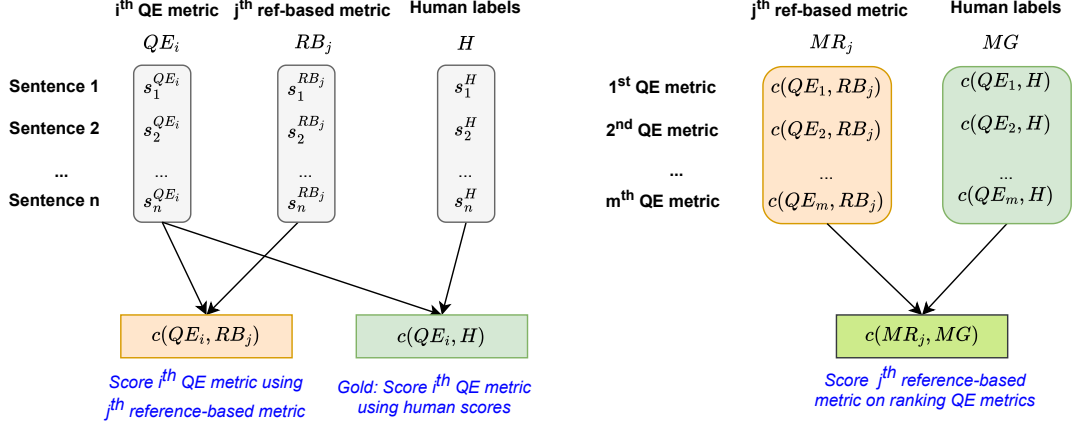


Figure 1: Illustration of our automatic QE evaluation approach.

where argmin^k returns the indices of k smallest elements and L_2 is the Euclidean distance.

Derive QE metrics Given the retrieved k -nearest neighbors, we derive QE metric s_j for each inference-time generated token y_j :

- k NN token distance: We calculate the average distance from y_j to its k -nearest tokens in the datastore:

$$s_j = \text{avg}_{(m,i) \in N_j} (L_2(d_j, \hat{d}_i^m)) \quad (4)$$

We assume the lower the distance, the better the translation quality, since the generated token is familiar to the MT model.

- k NN sentence similarity: We calculate the average cosine similarity between the whole inference-time generated sentence and the K sentences in the training data to which the k NN tokens belong:

$$s_j = \text{avg}_{(m,i) \in N_j} (\text{cos_sim}(\text{emb}(Y), \text{emb}(\hat{Y}^m))) \quad (5)$$

where cos_sim is the cosine similarity function, emb is the sentence embedding function. For sentence embedding, we use an external model instead of the MT model itself, since the external model won't be affected by artifacts in the MT training data.

We assume the higher the similarity, the better the translation quality, since the generated sentence is familiar to the MT model.

- Number of different k NN tokens: We count the number of distinct tokens amongst the retrieved k NN tokens:

$$s_j = |\{\hat{y}_i^m | (m, i) \in N_j\}| \quad (6)$$

We assume the higher the number, the lower the translation quality, since it means the neighbor cluster is not representative of any specific token, indicating that the model is uncertain about the generated representation.

- Model prediction equals retrieved k NN tokens: We count the number of retrieved k NN tokens that are the same as the model output token y_j :

$$s_j = |\{\hat{y}_i^m | (m, i) \in N_j \wedge \hat{y}_i^m = y_j\}| \quad (7)$$

We assume the higher the number, the better the translation, since it is easy for the model to map the representation to one single token.

Using these metrics, we can get quality scores on the token level. To get scores on the segment level, we take the average of the scores of tokens:

$$s_Y = \text{avg}_{j \in 1..|Y|} (s_j) \quad (8)$$

4 Automatic evaluation for Quality Estimation

Motivation Normally, to evaluate QE metrics, people calculate the correlation between QE-generated quality scores and human-generated quality scores on some MT output (Zerva et al., 2022). However, as discussed above, model-specific QE approaches such as k NN-QE provide quality scores on their own MT output, thus we cannot evaluate them using the available human quality scores on different MT outputs in the public benchmarks. Collecting human-generated quality scores again on this specific MT output would be costly in terms of time and human resources.

Therefore, we propose an automatic approach using reference-based metrics as gold standard to evaluate QE metrics.

Reference-based metrics as gold Recall that Quality Estimation takes only the source sentence and the MT translation for outputting a quality score, while reference-based metrics also make use of human gold-standard translation. As a result, reference-based metrics are usually more robust than QE metrics (Freitag et al., 2023). Therefore, we attempt to perform automatic evaluation for QE by calculating the correlation between the QE scores and the ref-based metrics scores. In other words, we are using the ref-based metrics scores in place of human-provided scores as the gold standard. We investigate scores at the segment level. Using this approach, we can flexibly generate gold-standard quality scores for any MT output, rather than relying on fixed human quality scores on some pre-made MT output.

Boosting reference-based metrics’ reliability Intuitively, it is important to have a robust reference-based metric since we are using it as gold standard for reference-free QE. One potential way to have more robust reference-based metrics is to increase the number of references. Therefore, we propose to use test datasets with multiple references, and additionally use a paraphraser tool to generate synthetic references.

Choosing reference-based metric We investigate whether reference-based metrics are good enough for evaluating QE metrics, and which reference-based metric is best suited. We gather different QE metric submissions on public shared tasks, and measure the correlation between the QE ranking created by human annotations and the QE ranking created by reference-based metrics. An illustration of the process is shown in Figure 1.

Specifically, assume we have n MT output segments, m QE metrics and p reference-based metrics. Let $QE_i, RB_j, H \in \mathbb{R}^{n \times 1}$ be the quality scores assigned to the MT translations by the i^{th} QE metric, the j^{th} reference-based metric and the human annotator respectively. The gold evaluation for the QE metrics is then:

$$MG = (c(QE_1, H), c(QE_2, H), \dots, c(QE_m, H)) \quad (9)$$

where c is a correlation function such as Spearman. The automatic evaluation for the QE metrics using

the j^{th} reference-based metric is:

$$MR_j = (c(QE_1, RB_j), c(QE_2, RB_j), \dots, c(QE_m, RB_j)) \quad (10)$$

The performance of the j^{th} reference-based metric on ranking QE metrics is then:

$$c(MR_j, MG) \quad (11)$$

Note that this is not the same as the performance of the j^{th} reference-based metric on scoring segment-level MT, which is defined as:

$$c(RB_j, H) \quad (12)$$

5 Experimental Setup

5.1 Automatic evaluation for Quality Estimation

Dataset In our experiments, we use the English – German data from two shared tasks: *WMT22 Quality Estimation* (Zerva et al., 2022) and *WMT22 QE as a Metrics* (Freitag et al., 2022). The *WMT22 Quality Estimation* shared task, which we refer to as *QE Task*, is specialized in evaluating Quality Estimation. The *WMT22 QE as a Metrics* shared task, which we refer to as *QE-M Task*, is meant for comparing QE metrics to reference-based metrics. Both shared tasks contain submissions from different QE systems, which is useful for us to investigate whether we can automatically rank these QE systems. Specifically, the QE-M Task data includes source sentences, reference sentences and translation sentences from multiple different MT systems from the WMT22 General MT task (Kocmi et al., 2022), along with human-labeled MQM quality score (Lommel et al., 2014) and QE submission scores on each translation sentence. The data from the QE Task is similar, except that (1) they only use data from the News domain rather than the full test set from the WMT22 General MT task (including the Conversation, Ecommerce, News and Social domains) and (2) the MT output is from a single MT system. More details can be found in Table 1.

Models and Tools We use Spearman’s rank correlation coefficient ρ for the calculation of automatic QE evaluation, i.e., the correlation function used in Equations 9, 10, 11 and 12. For creating synthetic references, we use a German paraphraser available on Huggingface[‡]. We consider

[‡]<https://huggingface.co/Lelon/t5-german-paraphraser-large>

	QE Task	QE-M Task
Domain	News	Multiple *
# sentence pairs	511	2,037
# references	2	2
# MT systems	1	14
# QE metrics	10	10

* Conversation, Ecommerce, News, Social

Table 1: Statistics of WMT22 Tasks on English–German.

different reference-based metrics to see which one is suitable for automatic QE evaluation, which includes: (1) lexical-based metrics: BLEU (Papineni et al., 2002), TER (Snover et al., 2006) and chrF (Popović, 2015); (2) embedding-based metrics: BERTScore (Zhang et al., 2019) and (3) neural-based metrics: BLEURT (Sellam et al., 2020), UniTE-MUP (Wan et al., 2022), COMET 22 (Rei et al., 2022a), xCOMET XL (Guerreiro et al., 2023) and MetricX-23 XL (Juraska et al., 2023).

5.2 Quality Estimation with k NN

Dataset We use the TED talks English–German bitext data from the evaluation campaign IWSLT 2014 (Cettolo et al., 2014) for training/fine-tuning MT models. The dataset includes 174,443 training sentences, 2,052 validation sentences and 4,698 testing sentences. For evaluation, we use TED test split, and additionally an out-of-domain test set from the WMT22 General task, i.e., the same data as for the automatic QE evaluation experiments.

For generating the train datastore, we use the train split of TED, i.e., the training data of the MT models. Additionally, we try to use an external, non-train datastore generated using the Europarl dataset (Koehn, 2005). From Europarl, we selected a subset with similar size as the TED training set to rule out the data size factor when comparing the external datastore to the TED datastore. We generate this external datastore similarly to the train datastore, where we perform inference with reference-forced decoding on Europarl using the TED-trained MT models.

MT models We consider two MT models: a model trained from scratch on TED, and a pre-trained DeltaLM model (Ma et al., 2021) fine-tuned on TED. The model trained from scratch uses the transformer base architecture from the Fairseq library (Ott et al., 2019), with 6 encoder layers, 6 decoder layers and embedding size of 512. Its vocabulary size is 10,112. The fine-

tuned DeltaLM model uses the DeltaLM base architecture with 12 encoder layers, 6 decoder layers and embedding size of 768. Its vocabulary size is 250,001. For the fine-tuned DeltaLM model, we build the datastore using only the fine-tuning data (TED), not the whole pretraining data of DeltaLM.

Automatic QE evaluation We focus on evaluating k NN-QE on the segment level. We use our automatic evaluation method described in Section 4. We calculate the Spearman correlation between segment-level scores generated by the QE metrics and gold scores generated by the reference-based MetricX-23 (Juraska et al., 2023), since we find MetricX-23 to be the most robust in QE ranking. Information about the token-level experiments on k NN-QE can be found in Appendix A.

Baselines We use the probability output from the MT model as an unsupervised QE baseline. We take the average of the probability for each token to get the segment-level score. The higher the probability, the better the quality, as it is an indication that the MT model is confident. Since our k NN-QE approach is unsupervised, we choose the supervised WMT22 COMET-Kiwi model (Rei et al., 2022b) as an upper-bound for the performance.

Tools For training/fine-tuning MT models, we use the Fairseq library (Ott et al., 2019). For generating the datastore and retrieving k NN samples, we use the k NN-box toolkit (Zhu et al., 2023), which makes use of Faiss (Johnson et al., 2019) for efficient similarity search. For embedding sentences, we use an external model from Huggingface[§]. Experiments were conducted on an Nvidia TITAN RTX GPU with 25 GB of memory.

6 Results and Discussion

6.1 Automatic QE evaluation

Overall performance The performance of difference reference-based metrics on ranking QE submissions on the two WMT22 shared tasks is shown in Table 2. MetricX-23 XL performs the best on the QE-M Task with 0.939 Spearman correlation to human-based ranking. BLEU performs the best on ranking QE metrics on the QE Task with 0.721 Spearman correlation to human-based ranking. Given these high correlations, we con-

[§]<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

clude that using reference-based metrics is sufficient for automatically evaluating QE metrics.

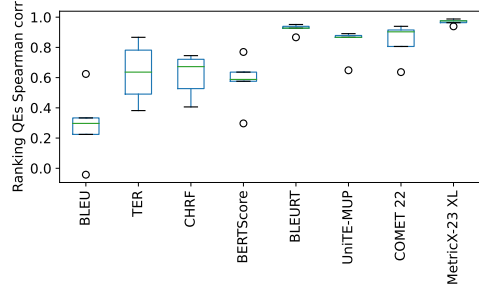
	QE Task	QE-M Task
BLEU	0.721	0.333
TER	0.685	0.782
chrF	0.564	0.745
BERTScore	0.636	0.576
BLEURT	0.442	0.927
UniTE-MUP	0.321	0.867
COMET 22	0.273	0.903
xCOMET XL *	0.358	-
MetricX-23 XL	0.261	0.939

*: xCOMET models are trained on WMT22 data except for the News domain, thus only valid to be tested on the QE Task data.

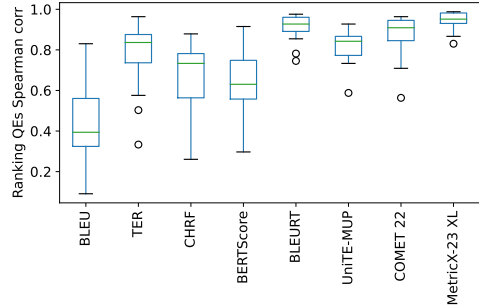
Table 2: Overall performance of different reference-based metrics on ranking QE on two public shared tasks.

We take a closer look at the correlations on the QE Task. It is quite surprising that BLEU has the highest correlation in QE ranking, since BLEU has recently been shown to have worse evaluation performance than other neural-based metrics (Freitag et al., 2022). However, BLEU’s ranking correlation is quite low on the QE-M Task data as expected. A similar pattern can be observed where MetricX-23 XL has unexpectedly low performance on the QE Task, but good performance on the QE-M task. We assume that the unexpected ranking performance of the metrics on the QE Task data is due to the narrow scope of QE Task: it considers the output of a single MT model on a single domain. On the other hand, the QE-M Task data is on multiple MT systems output on multiple domains. Therefore, we suspect that the results on the QE-M Task are potentially more generalizable, and that MetricX-23 XL is the best metric for ranking QE metrics. Our following experiment results show evidence that supports this assumption.

MetricX-23 XL robustness We collect the performance of reference-based metrics on QE ranking across different domains and different MT systems’ output, as can be seen in Figure 2a and Figure 2b, respectively. Generally, the neural-based metrics have better performance than the lexical-based and embedding-based metrics. Their scores are higher and more consistent across different domains and MT systems’ output. Among the neural-based metrics, MetricX-23 XL has the best performance in terms of score and consistency.



(a) Group by domains in WMT 22 General.



(b) Group by MT systems participated in WMT 22 General.

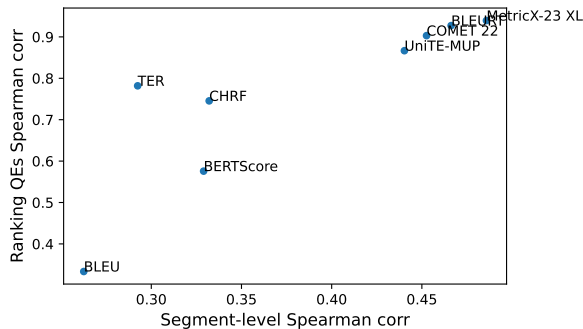
Figure 2: QE ranking performance across different factors.

Evaluating segments versus evaluating QE metrics

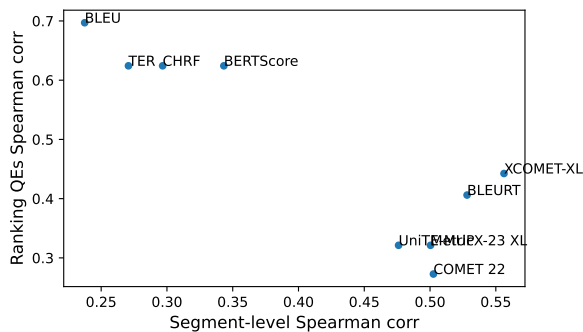
We investigate whether better performance on evaluating MT segments (Equation 12) means better performance on evaluating QE metrics (Equation 11) for reference-based metrics. Figure 3 shows that this is not always the case. For example, in Figure 3a on the QE-M Task data, TER and chrF have low performance on segment-level evaluation, but have decent performance on ranking QEs. However, both of them are still worse than MetricX-23 XL. In Figure 3b on the QE Task data, the pattern is even more unexpected, where the lexical-based metrics have significantly better performance in QE ranking than the neural-based metrics, while being worse at evaluating segment-level MT output. However, we suspect that this is due to the QE Task data being specific on a single domain and a single MT system’s output, thus the result is not representative. The following experiment result supports this assumption.

Importance of a broad-ranged test set

We perform the same experiment on segment-level evaluation performance versus QE ranking performance on the QE-M Task, but limit it to a single domain and a single MT system. In Figure 4a, on a single MT system output on all domains, the neural-based metrics perform well on both segment-level evaluation and QE ranking as expected. However,



(a) QE-M Task: multiple domains, multiple MT systems.



(b) QE Task: single domain, single MT system.

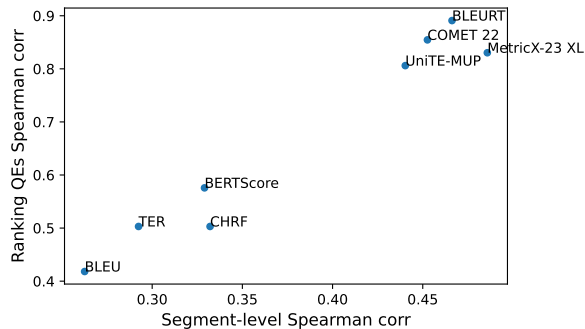
Figure 3: Correlation between the performance on evaluating translation segments and the performance on QE ranking.

on a single MT system output on a single domain (Figure 4b), neural-based metrics have worse QE ranking performance than some lexical-based metrics, while still doing well on segment-level evaluation. It can be concluded that reference-based metrics can have unexpected performance when the testing condition is too narrow. Therefore, it is important to perform evaluation on a broad-ranged test set with multiple domains so that we can rely on neural reference-based metrics for QE ranking.

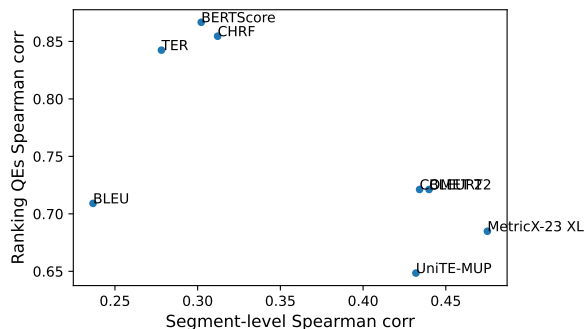
Importance of references: quantity and quality Figure 5 shows the effect of references on reference-based metrics’ performance on QE ranking. Having two human-created references is better than one, showing that increasing the quantity of references helps improve performance. However, adding synthetic references created by paraphrasing decreases the performance to some extent. This shows that it is important to add high-quality references, otherwise it might have the opposite effect of harming the overall performance.

6.2 Quality Estimation with k NN

We report on the segment-level performance of k NN-QE. Experiments on k NN-QE performance on the token level can be found in Appendix A.



(a) QE-M Task: all domains, single MT system (MT system: comet_bestmbrMT).



(b) QE-M Task: single domain (Social), single MT system (comet_bestmbrMT).

Figure 4: Correlation between the performance on evaluating translation segments and the performance on QE ranking, limited by MT system and domain.

k NN-QE better than MT probability, but worse than supervised QE The performance of our k NN-QE approach is shown in Table 3. We consistently observe that the performance of the k NN token distance metric is better than the other k NN-QE metrics. k NN token distance (Row 3) has better performance than the probability baseline (Row 1) by 0.1 increase in Spearman correlation to humans in most cases. However, it still falls behind the supervised QE baseline (Row 2). Ensembling all four k NN-QE metrics gives improvement in performance, but not significant, as it is only \approx

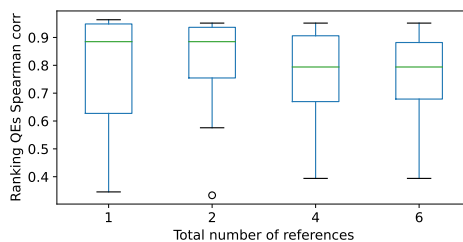


Figure 5: QE-M Task: Effect of number of references. The first 2 boxes use human references only, while the last 2 boxes also include synthetic references created by paraphrasing.

			Transformer Scratch		Fine-tuned DeltaLM	
			TED	WMT22	TED	WMT22
Baselines						
1	Probability		0.535	0.525	0.462	0.423
2	Supervised QE		0.773	0.793	0.705	0.771
<i>k</i>NN-QE						
3	TED	<i>k</i> NN token distance ^a	0.650	0.623	0.575	0.438
4	datastore	<i>k</i> NN sentence similarity ^a	0.570	0.553	0.527	0.398
5		<i>k</i> NN nr. distinct tokens ^b	0.475	0.469	0.423	0.336
6		<i>k</i> NN tokens = output token ^b	0.489	0.497	0.410	0.348
7		Ensemble ^c	0.652	0.627	0.576	0.439
8	20% TED	<i>k</i> NN token distance ^a	0.620	0.601	0.554	0.412
9	datastore	<i>k</i> NN sentence similarity ^a	0.532	0.486	0.496	0.373
10		<i>k</i> NN nr. distinct tokens ^b	0.498	0.494	0.407	0.373
11		<i>k</i> NN tokens = output token ^b	0.491	0.507	0.390	0.353
12		Ensemble ^c	0.622	0.604	0.555	0.413
13	Europarl	<i>k</i> NN token distance ^a	0.546	0.514	0.543	0.414
14	datastore	<i>k</i> NN sentence similarity ^a	0.121	0.246	0.103	0.051
15	(≠ train)	<i>k</i> NN nr. distinct tokens ^b	0.383	0.351	0.320	0.271
16		<i>k</i> NN tokens = output token ^b	0.437	0.465	0.384	0.335
17		Ensemble ^c	0.548	0.517	0.544	0.415

^a: Number of neighbors $k = 1$.

^b: Number of neighbors $k = 10$.

^c: Ensembling from the other four KNN-QE metrics.

Table 3: Overall performance of *k*NN-QE on the segment level.

0.002 points higher than the performance of the *k*NN token distance metric alone.

Performance diminishes with fine-tuned MT on out-of-domain test set From Table 3, we can see the performance change when moving from the in-domain test set (TED) to out-of-domain test set (WMT22). For the Transformer MT model trained from scratch on TED (“Transformer Scratch”), our approach works for both in-domain and out-of-domain test sets, where it outperforms the probability baseline. On the other hand, for the DeltaLM model fine-tuned on TED, our approach only outperform the probability baseline on the in-domain test set. On the out-of-domain test set, it performs similar or worse than the probability baseline. This is possibly due to the fine-tuned DeltaLM model being pretrained on other data than TED, thus having more knowledge on the out-of-domain test set which is not identifiable if we only use the datastore on TED. It can be concluded that (1) *k*NN-QE works best if we build the datastore using the training data of the model, not only fine-tuning data and

(2) with the appropriate training datastore, the approach works for out-of-domain test sets.

Reducing the datastore has less negative effect than expected As can be seen in Figure 6, generally, the QE performance of the *k*NN token distance and *k*NN sentence similarity metrics increases as the portion of training data used to create the datastore increases. However, the QE performance only increases drastically if we increase the datastore until 20%, afterward, it starts to flatten. For the other two metrics, the QE performance slightly fluctuates with different datastore sizes. More detailed numbers can be seen in Table 3, where the QE performance using 20% datastore is only worse than using the full TED datastore by ≈ 0.03 reduction in Spearman correlation. This is a positive observation, since building a smaller datastore would be more memory-efficient and inference-speed-efficient.

Switching to non-train datastore hurts performance As can be seen from Table 3, changing to the Europarl datastore reduces the performance of

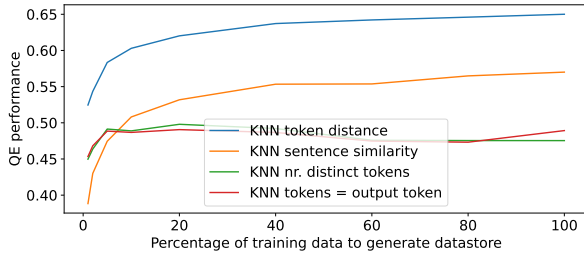


Figure 6: Effect of reduced datastore. Experiment conducted with Transformer Scratch MT model on TED. Similar patterns observed with fine-tuned DeltaLM and WMT 22 data.

k NN-QE. This is expected, since a non-train datastore would not be representative of the MT model’s knowledge. However, using the Europarl datastore for k NN-QE still works to some extent, as the QE correlation to humans is still quite high, at around 0.5 using the k NN token distance metrics (Row 13). This is potentially due to the use of forced decoding on reference: in the datastore, we use reference translation as prefix to generate each token, making the generated tokens have higher quality. Thus the closer the inference-time generated tokens are to the high-quality ones in the datastore, the more likely that they also have high quality.

Interestingly, using this same-size but non-train datastore leads to worse performance than using only 20% of the train datastore, which further strengthens the importance of having a datastore that represents the model’s knowledge.

We also observe that the negative impact of switching to a non-train datastore is less significant for the fine-tuned DeltaLM model than the Transformer model trained from scratch. This is potentially due to DeltaLM’s pretraining data containing the same or similar data to the Europarl data, thus the Europarl datastore represents the knowledge of the fine-tuned DeltaLM model to some extent.

Effect of number of neighbors As can be seen in Figure 7, k NN token distance and k NN sentence similarity metrics only need a small number of retrieved neighbors. Their performance decreases as the number of nearest neighbors increases. This is an indication that only the distance of the inference-time generated token to its closest training neighbor matters for these two metrics. However, for the other 2 metrics, i.e., number of distinct k NN tokens and number of k NN tokens same as model output, the higher the number of neighbors retrieved the better. This is due

to these two metrics only comparing the surface-level token output, thus retrieving a small number of neighbors doesn’t provide as much information. Based on these observations, we choose the number of neighbors to be $k = 1$ for k NN token distance and k NN sentence similarity metrics and $k = 10$ for the other two metrics to report in the main Table 3.

Observe that with different numbers of nearest neighbors, the k NN token distance metric still performs the best. This means that we can go for this metric in practice with a small number of retrieved neighbors, which benefits the inference speed. Combining the small value of $k = 1$ with the reduced 20% TED datastore, we observe around 19% increase in inference time when applying k NN-QE to the generation process.

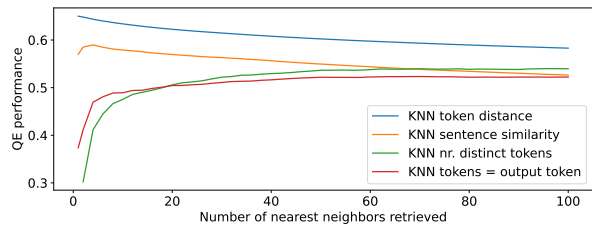


Figure 7: Effect of number of neighbors. Experiment conducted with Transformer Scratch MT on TED. Similar patterns observed with fine-tuned DeltaLM and WMT 22 data.

7 Conclusion

In this paper, we proposed k NN-QE – a model-specific, unsupervised Quality Estimation approach which exploits the information from the MT model’s training data. We also propose an automatic QE evaluation method for such model-specific QE approaches, which make use of reference-based metrics. Our experiments show that this automatic evaluation method is sufficient, and that the reference-based MetricX-23 XL is the most suitable. Using this automatic QE evaluation method, we found that k NN-QE performs better than the MT probability baseline, but still falls behind the supervised QE approach. We also find that our approach works with a small number of retrieved neighbors and a small portion of the training datastore, making it more memory- and time-efficient to be used in practice. For future work, we can explore whether this method is applicable to other types of generative models, such as the currently prominent Large Language Models.

Acknowledgements

This work is supported by the Helmholtz Programme-oriented Funding, with project number 46.24.01, project name AI for Language Technologies. It also received partial support from the European Union’s Horizon research and innovation programme under grant agreement No 101135798, project Meetween (My Personal AI Mediator for Virtual MEETings BetWEEN People). It was partly performed on the HoreKa supercomputer funded by the Ministry of Science, Research and the Arts Baden-Württemberg and by the Federal Ministry of Education and Research.

References

- Blain, Frederic, Chrysoula Zerva, Ricardo Rei, Nuno M. Guerreiro, Diptesh Kanojia, José G. C. de Souza, Beatriz Silva, Tânia Vaz, Yan Jingxuan, Fatemeh Azadi, Constantin Orasan, and André Martins. 2023. Findings of the WMT 2023 shared task on quality estimation. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 629–653, Singapore, December. Association for Computational Linguistics.
- Cettolo, Mauro, Jan Niehues, Sebastian Stüker, Luisa Bentivogli, and Marcello Federico. 2014. Report on the 11th IWSLT evaluation campaign. In Federico, Marcello, Sebastian Stüker, and François Yvon, editors, *Proceedings of the 11th International Workshop on Spoken Language Translation: Evaluation Campaign*, pages 2–17, Lake Tahoe, California, December 4-5.
- Fomicheva, Marina, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. 2020. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:539–555.
- Freitag, Markus, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Freitag, Markus, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore, December. Association for Computational Linguistics.
- Guerreiro, Nuno M, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2023. xcomet: Transparent machine translation evaluation through fine-grained error detection. *arXiv preprint arXiv:2310.10482*.
- Johnson, Jeff, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.
- Juraska, Juraj, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore, December. Association for Computational Linguistics.
- Khandelwal, Urvashi, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2019. Generalization through memorization: Nearest neighbor language models. *arXiv preprint arXiv:1911.00172*.
- Khandelwal, Urvashi, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Nearest neighbor machine translation. *arXiv preprint arXiv:2010.00710*.
- Kocmi, Tom, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United

- Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Koehn, Philipp. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of Machine Translation Summit X: Papers*, pages 79–86, Phuket, Thailand, September 13–15.
- Lommel, Arle, Aljoscha Burchardt, Maja Popović, Kim Harris, Eleftherios Avramidis, and Hans Uszkoreit. 2014. Using a new analytic measure for the annotation and analysis of MT errors on real data. In Cettolo, Mauro, Marcello Federico, Lucia Specia, and Andy Way, editors, *Proceedings of the 17th Annual Conference of the European Association for Machine Translation*, pages 165–172, Dubrovnik, Croatia, June 16–18. European Association for Machine Translation.
- Lu, Yu, Jiali Zeng, Jiajun Zhang, Shuangzhi Wu, and Mu Li. 2022. Learning confidence for transformer-based neural machine translation. In Muresan, Smaranda, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2353–2364, Dublin, Ireland, May. Association for Computational Linguistics.
- Ma, Shuming, Li Dong, Shaohan Huang, Dongdong Zhang, Alexandre Muzio, Saksham Singhal, Hany Hassan Awadalla, Xia Song, and Furu Wei. 2021. DeltaLM: Encoder-decoder pre-training for language generation and translation by augmenting pretrained multilingual encoders.
- Niehues, Jan and Ngoc-Quan Pham. 2019. Modeling confidence in sequence-to-sequence models. In van Deemter, Kees, Chenghua Lin, and Hiroya Takamura, editors, *Proceedings of the 12th International Conference on Natural Language Generation*, pages 575–583, Tokyo, Japan, October–November. Association for Computational Linguistics.
- Ott, Myle, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Isabelle, Pierre, Eugene Charniak, and Dekang Lin, editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Popović, Maja. 2015. chrF: character n-gram F-score for automatic MT evaluation. In Bojar, Ondřej, Rajan Chatterjee, Christian Federmann, Barry Haddow, Chris Hokamp, Matthias Huck, Varvara Logacheva, and Pavel Pecina, editors, *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal, September. Association for Computational Linguistics.
- Rei, Ricardo, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Rei, Ricardo, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Rikters, Matīss and Mark Fishel. 2017. Confidence through attention. In Kurohashi, Sadao and Pascale Fung, editors, *Proceedings of Machine Translation Summit XVI: Research Track*, pages 299–311, Nagoya Japan, September 18 – September 22.
- Sellam, Thibault, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In Jurafsky, Dan, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online, July. Association for Computational Linguistics.
- Snover, Matthew, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge,

- Massachusetts, USA, August 8-12. Association for Machine Translation in the Americas.
- Vieira, Lucas Nunes, Minako O’Hagan, and Carol O’Sullivan. 2021. Understanding the societal impacts of machine translation: a critical review of the literature on medical and legal use cases. *Information, Communication & Society*, 24(11):1515–1532.
- Wan, Yu, Dayiheng Liu, Baosong Yang, Haibo Zhang, Boxing Chen, Derek Wong, and Lidia Chao. 2022. UniTE: Unified translation evaluation. In Muresan, Smaranda, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8117–8127, Dublin, Ireland, May. Association for Computational Linguistics.
- Wang, Ke, Yangbin Shi, Jiayi Wang, Yuqi Zhang, Yu Zhao, and Xiaolin Zheng. 2021. Beyond glass-box features: Uncertainty quantification enhanced quality estimation for neural machine translation. In Moens, Marie-Francine, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4687–4698, Punta Cana, Dominican Republic, November. Association for Computational Linguistics.
- Yankovskaya, Elizaveta, Andre Tättar, and Mark Fishel. 2018. Quality estimation with force-decoded attention and cross-lingual embeddings. In Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névéal, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 816–821, Belgium, Brussels, October. Association for Computational Linguistics.
- Zerva, Chrysoula, Frédéric Blain, Ricardo Rei, Piyawat Lertvittayakumjorn, José G. C. de Souza, Steffen Eger, Diptesh Kanojia, Duarte Alves, Constantin Orăsan, Marina Fomicheva, André F. T. Martins, and Lucia Specia. 2022. Findings of the WMT 2022 shared task on quality estimation. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névéal, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 69–99, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Zhang, Tianyi, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Zhang, Pei, Baosong Yang, Hao-Ran Wei, Dayiheng Liu, Kai Fan, Luo Si, and Jun Xie. 2022. Competency-aware neural machine translation: Can machine translation know its own translation quality? In Goldberg, Yoav, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4959–4970, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.
- Zhu, Wenhao, Qianfeng Zhao, Yunzhe Lv, Shujian Huang, Siheng Zhao, Sizhe Liu, and Jiajun Chen. 2023. knn-box: A unified framework for nearest neighbor generation.

Test set	QE method	Pearson Correlation	F1-score
TED	Probability	0.07	0.08
	k NN token distance ^a	0.21	0.27
News	Probability	0.14	0.12
	k NN token distance ^a	0.21	0.36

^a: Number of neighbors $k = 1$.

Table 4: k NN-QE performance on the token level.

Appendix A. k NN-QE on token level

A.1 Experimental Setup

Motivation for manual evaluation To evaluate the performance of k NN-QE on the token level, we need gold-standard token-level quality labels. On the segment level, we have proposed an automatic QE evaluation method (Section 4) by using segment-level quality scores made by reference-based metrics as gold standard instead of human quality scores. In principle, we can do the same for token-level evaluation, by finding a reference-based metric that provides quality labels on the token level.

However, the performance of reference-based metrics on the token level is usually not as good as on the segment level. For example, the xCOMET metric provides both error-span prediction and segment-level quality scores. Their segment-level quality scores correlate well with human MQM scores, at 0.653 Pearson. Meanwhile, the error-span prediction performance is quite poor, at 0.320 F1 score (although they are still very useful when being aggregated to provide segment scores). This is reasonable, since more fine-grain evaluation tends to be more difficult.

Due to the not-yet-perfect token-level performance of reference-based metrics, we choose not to use them as gold standard to evaluate k NN-QE. We instead opt for performing manual annotation on the token level of the MT output to evaluate k NN-QE.

Manually annotated data for evaluation We manually annotated the MT output on the token level. We annotate 2110 tokens from 100 output sentences on TED data (in-domain test set) and 3503 tokens from 100 output sentences on the News test data (out-of-domain test set).

Metrics Recall that for each generated subword, k NN-QE provides a quality score. To report the performance of k NN-QE on the token level, we

use two metrics: Pearson Correlation and F1-score. For Pearson Correlation, we treat the human-annotated labels as continuous scores (0 representing a *BAD* token, 1 representing a *OK* token), and calculate its correlation to the k NN-QE scores. For the F1-score, we turn the continuous k NN-QE scores into binary labels using a threshold. We choose a threshold that maximizes the F1-score.

Baseline We compare the performance of k NN-QE to an unsupervised baseline using probability output from the MT model.

Experiment scope Since it is difficult to perform manual evaluation on a large scale, we limit the scope of our experiment on the token level. In this experiment, we only report on the Transformer model trained from scratch on TED and our best k NN-QE metric, i.e., k NN token distance. Due to this small scale, we only include the token-level experiment here in the Appendix for more information, rather than including it in the main part of the paper.

A.2 Results and Discussion

As can be seen from Table 4, our k NN-QE outperforms the MT probability baseline. This is an indication that k NN-QE also works on the token level. Additionally, we observe that the QE performance is generally better on the out-of-domain test set.

SubMerge: Merging Equivalent Subword Tokenizations for Subword Regularized Models in Neural Machine Translation

Haiyue Song¹ Francois Meyer² Raj Dabre¹
Hideki Tanaka¹ Chenhui Chu³ Sadao Kurohashi^{3,4}

¹ NICT, Japan ² University of Cape Town, South Africa

³ Kyoto University, Japan ⁴ NII, Japan

{haiyue.song, raj.dabre, hideki.tanaka}@nict.go.jp,
francois.meyer@uct.ac.za,
{chu, kuro}@i.kyoto-u.ac.jp

Abstract

Subword regularized models leverage multiple subword tokenizations of one target sentence during training. Previous decoding algorithms select one tokenization during inference, leading to the underutilization of knowledge learned about multiple tokenizations. To address this, we propose the **SubMerge** algorithm to rescue the ignored **Sub**word tokenizations through **Merging** equivalent ones during inference. SubMerge is a nested search algorithm where the outer beam search treats words as the minimal units, and the inner beam search provides a list of word candidates and their probabilities by merging subword tokenizations that form the same word. Experimental results on six machine translation datasets show more accurate word probability estimation and higher translation quality using SubMerge than beam search. Additionally, we provide time complexity analysis and investigate the effect of different beam sizes, training set sizes, dropout rates, and whether it is effective on non-regularized models.

1 Introduction

Despite the end-to-end nature that makes neural machine translation (NMT) (Sutskever et al., 2014; Bahdanau et al., 2014; Vaswani et al., 2017; Gehring et al., 2017) the most prevalent and convenient approach for machine translation (MT), subword tokenization (Sennrich et al., 2016b;

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

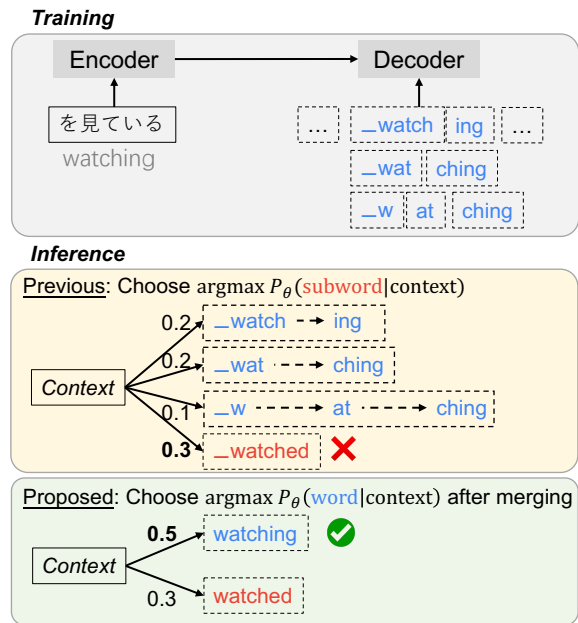


Figure 1: Subword regularized models suffer from discrepancies between training and inference, where they are trained on multiple target tokenizations and generate one. We propose to merge **equivalent subword tokenizations** that compose the same word with different conditional probabilities during the inference.

Provilkov et al., 2020; Kudo and Richardson, 2018; Kudo, 2018a) remains an indispensable pre-processing step for most NMT systems. Subword vocabularies address the out-of-vocabulary problem of word-based NMT systems (Kalchbrenner and Blunsom, 2013; Bahdanau et al., 2014; Luong et al., 2015) by reducing new words to known subwords, while avoiding the high computational cost of character-based NMT systems (Gupta et al., 2019; Kim et al., 2016; Costa-jussà and Fonollosa, 2016; Ling et al., 2015; Cherry et al., 2018) by enabling much shorter input and output sequences.

Deterministic segmenters like Byte-Pair Encoding (BPE) (Sennrich et al., 2016b) and Sentence-

Piece (Kudo, 2018a) are widely used due to their simplicity and effectiveness. They are *deterministic* in the sense that they consistently generate the same tokenization for a given sentence. NMT models trained on consistent subword tokenizations typically allocate the majority of a sentence’s true probability (considering all potential tokenizations by marginalizing over them) to its specific tokenization (Cao and Rimell, 2021), except for out-of-domain data (Chirkova et al., 2023). Therefore, the probability of the sentence approximately equals the probability of that tokenization.

On the other hand, stochastic segmenters such as subword regularization methods (Provilkov et al., 2020; Kudo, 2018a) produce multiple tokenizations of a given sentence during training, as illustrated in Figure 1. As a data augmentation method, models trained on regularized data usually outperform those trained on non-regularized data, especially in low-resource scenarios. However, this causes a discrepancy between training and inference. During training, the model learns to generate multiple target tokenizations for each source sentence and learns to distribute the probability of a target sentence across all the tokenizations. During inference, greedy or beam search approximates the *single* highest probability tokenization. This causes a discrepancy - the probability of a target tokenization diverges drastically from the probability of a target sentence. The inaccurate probability estimation of the next word during inference in turn leads to a degradation in translation quality. The way to overcome this is to incorporate the marginal likelihood of the next words during decoding for the subword regularized models.

To this end, we propose SubMerge, a decoding algorithm that aggregates probabilities from exponentially many tokenizations for a sentence by merging subword tokenizations that form the same word. The property of BPE-dropout (Provilkov et al., 2020) that each word is individually segmented makes aggregating probabilities from exponentially many tokenizations theoretically possible. As for the implementation, SubMerge is a nested beam search approach. In the outer beam search, we hide the detail of possible subword tokenizations of the word, treating words as minimal units. This ensures that the outer beam is unaware of and unaffected by the subword tokenizer. In the inner beam search, we limit the search space within the word boundary. The inner beam search

finds the n -best tokenizations, merges equivalent ones, and returns a list of words and the corresponding probabilities.

Previous attempts to estimate marginal likelihood over tokenizations include summing over n -best tokenizations (Cao and Rimell, 2021) and using importance sampling (Chirkova et al., 2023). However, these algorithms focus on perplexity estimation, assuming the output is already in hand. In our approach, we perform marginal likelihood estimation for the next words along with the inference process, aiming to improve not only the estimation precision but also the translation quality. In a nutshell, our contributions are as follows:

- We propose SubMerge, a nested beam search algorithm for generating text with subword regularized models. It merges equivalent subword tokenizations for the next words, thereby enhancing probability estimation precision and translation quality.
- Experimental results on six machine translation datasets demonstrate significant improvements in estimating the underlying word perplexity computation for a model and its translation quality.
- We provide analyses of time complexity, various beam sizes, the selection of the inner searching function, and the impact of hyperparameters.

2 Preliminaries

This section formulates the objective of the inference process of NMT models, highlights the distinction introduced by subword regularized models, and introduces how we address it.

Inference Objective An NMT model with parameters θ during inference is to obtain $\arg \max_Y P_\theta(Y|X)$ where X and Y are the source and target sentences in plain text form. For subword-based NMT models, we tokenize X into a sequence of tokens during both training and inference. We tokenize Y during the training and try to predict a sequence of tokens that compose Y during inference. We use two tokenizers $\tau_S(X) = \mathbf{x}$, where $\mathbf{x} = (x_1, \dots, x_n)$ and $\tau_T(Y) = \mathbf{y}$, where $\mathbf{y} = (y_1, \dots, y_m)$. Each subword x_i or y_i is a non-empty substring of the text X or Y in a finite-size subword vocabulary predefined by the source or

target tokenizer. In theory,

$$P_\theta(Y|X) \neq P_\theta(\mathbf{y}|\mathbf{x}), \quad (1)$$

because there are multiple tokenizations of X and Y (besides \mathbf{x} and \mathbf{y}) that the model P_θ would assign non-zero probabilities to (Cao and Rimell, 2021).

Non-regularized Models For NMT models using deterministic tokenizers such as BPE (Sennrich et al., 2016b), tokenization function $\tau(\cdot)$ is a bijective function, and we can approximate the objective using one tokenization with a small gap (less than 0.5%) (Chirkova et al., 2023):

$$P_\theta(Y|X) \approx P_\theta(\mathbf{y}|\mathbf{x}). \quad (2)$$

Therefore, we can use $\arg \max_{\mathbf{y}} P_\theta(\mathbf{y}|\mathbf{x})$ to approximate $\arg \max_Y P_\theta(Y|X)$ with greedy or beam search in inference. This allows us to identify the next tokens with high conditional probabilities without concern for the discrepancy between the probability of raw text Y and of the particular tokenization \mathbf{y} .

Subword Regularized Models For NMT models using stochastic tokenizers (Provilkov et al., 2020), the tokenization function τ yields multiple tokenizations for one sentence. That is $\tau_S(X) = \mathbf{x} \in \mathcal{V}_S(\mathcal{X})$ where $\mathbf{x} \sim P_{\tau_S}(\mathbf{x}|X)$. Similar for Y . In this case, the number of possible segmentation $\mathcal{V}_S(\mathcal{X})$ increases exponentially according to the length of X , which deviates $P_\theta(\mathbf{y}|\mathbf{x})$ drastically from $P_\theta(Y|X)$, thus it requires marginalization over all possible tokenizations:

$$P_\theta(Y|X) = \sum_{\mathbf{x} \in \mathcal{V}_S(X)} \sum_{\mathbf{y} \in \mathcal{V}_T(Y)} P_\theta(\mathbf{y}|\mathbf{x}) P_{\tau_S}(\mathbf{x}|X). \quad (3)$$

This study focuses on better estimating the marginal likelihood of the target side, so we simplify Eq. (3) by using the most probable source tokenization $\arg \max_{\mathbf{x} \in \mathcal{V}_S(X)} P_{\tau_S}(\mathbf{x}|X)$ and remove the effect of the source tokenizer, resulting in:

$$P_\theta(Y|X) \approx \sum_{\mathbf{y} \in \mathcal{V}_T(Y)} P_\theta(\mathbf{y}|\mathbf{x}). \quad (4)$$

Inference for Subword Regularized Models We propose SubMerge to approximate Eq. (4) by introducing an intermediate variable, word tokenizations $\mathbf{w} = (w_1, \dots, w_n)$, generated by a word to-

kenizer $\tau_W(\cdot)$ which is a bijective function.¹ The problem is simplified as:

$$P_\theta(Y|\mathbf{x}) = P_\theta(\mathbf{w}|\mathbf{x}) = \prod_{i=1}^n P_\theta(w_i|\mathbf{w}_{<i}, \mathbf{x}). \quad (5)$$

We estimate $P_\theta(w_i|\mathbf{w}_{<i}, \mathbf{x})$ by summing over probabilities of subword tokenizations for one word w_i where the search space is much smaller compared to the search space of tokenizations of a whole sentence in Eq. (4):

$$P_\theta(w_i|\mathbf{w}_{<i}, \mathbf{x}) \approx \sum_{\mathbf{y}' \in \mathcal{V}_T(w_i)} P_\theta(\mathbf{y}'|\mathbf{w}_{<i}, \mathbf{x}). \quad (6)$$

In practice, since the decoder only takes subword as input, we feed the best subword tokenization of the next word w_i , which is $\arg \max_{\mathbf{y}' \in \mathcal{V}_T(w_i)} P_\theta(\mathbf{y}'|\mathbf{w}_{<i})$.

In this way, the probability of the target sentence is accurately calculated through a deterministic word tokenization as shown in Eq. (5), where the probability estimation of each word is precisely estimated through marginal likelihood estimation shown in Eq. (6). We implement Eq. (5) with the outer beam search as introduced in Section 3.2 and Eq. (6) with our inner beam search as introduced in Section 3.3.

3 Methodology

3.1 Overview of SubMerge

An overview of the SubMerge algorithm is shown in Figure 2. It is a nested beam search decoding algorithm that contains an outer beam search as explained in Section 3.2 and an inner beam search with subword merging post-processing as explained in Section 3.3. The outer beam search selects from a list of words considering the conditional probability in each step and estimates the most probable sentence $\arg \max_Y P_\theta(Y|X)$. The inner beam estimates the conditional probability of words in Eq. (6) by merging the probabilities of different subword tokenizations of the same words.

3.2 Outer Beam Search

The outer beam search algorithm is shown in Algorithm 1. It follows the standard beam search approach, where the difference is that words serve

¹That is $\tau_W(Y) = \mathbf{w}$. Note that word tokenizer is not a bijective function for languages such as Japanese or Chinese. For these languages, we can use specific word segmenters such as Jumanpp or Stanford Word Segmenter, which are bijective.

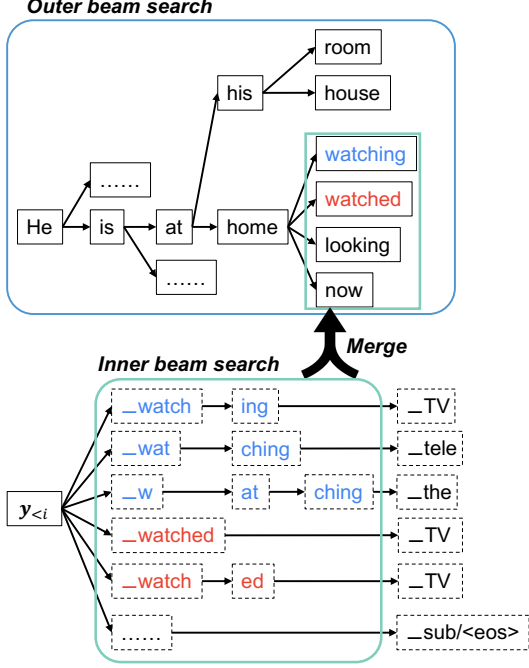


Figure 2: Overview of SubMerge. It contains an outer beam search that views words as minimal units. The candidate words and their probabilities are obtained from merging subword tokenizations in the n -best list of the inner beam search.

as the basic units. At each time step t in line 3, we explore each state s in the queue B_{t-1} that saves the best results in the previous step. When s is not finished, the candidates of the next words are obtained from a call to the inner beam search shown in line 9. Each state in the outer beam search queue contains the probability of the generation, the previous words, and their most probable tokens. Each state s' from the inner beam search contains the probability of the possible next word, the next word itself, and the most probable subword tokenization of that word. We add the new state to B_t shown in line 14 and only save the top- K ones shown in line 16. The most probable subword tokenization is used as the contextual input in the next decoding step.

In practice, we take the logarithm ($\log(\cdot)$) of the probabilities for computational precision. We implemented early stopping after all sequences reach the special end-of-sentence ($\langle eos \rangle$) token.

3.3 Inner Beam Search

The inner beam search is shown in Algorithm 2. It consists of two parts: 1) a token-level beam search within the word boundary and 2) post-processing to merge probabilities from equivalent subword tokenizations that compose the same word.

The first part is similar to that of the outer beam

Algorithm 1: OuterBeamSearch

Data: Beam width K , max length T
Result: Best sequence of states

```

1 Initialization:
2  $B_0 \leftarrow \{(0, [], [])\}$ ;
3 for  $t \leftarrow 1$  to  $T$  do
4    $B_t \leftarrow \emptyset$ ;
5   foreach  $s \in B_{t-1}$  do
6     if  $s$  reaches  $\langle eos \rangle$  then
7        $B_t.append(s)$ ;
8       continue;
9     foreach  $s' \in InnerBS(s[2])$  do
10       $score, word, toks = s'$ ;
11       $score \leftarrow s[0] + score$ ;
12       $words \leftarrow s[1] + words$ ;
13       $toks \leftarrow s[2] + toks$ ;
14       $B_t.append((score, words, toks))$ ;
15   Sort  $B_t$  by scores in descending order;
16    $B_t \leftarrow B_t[:K]$ 
17 return  $B_T$ 

```

search. The stopping criteria of one sequence are reaching the start of the next word (with the start-of-word indicator ' _ ' Unicode U+2581) or the $\langle eos \rangle$ token, where this stopping token will not be added to the token list. Otherwise, the exploration of the sequence continues according to the next subword probability distribution given by the decoder. During the post-processing part, we remove special tokens and spaces during the detokenization of a token list to form the word and return a list of words with their probabilities. The time complexity of SubMerge is $O(T \cdot K^3)$, where T is the sentence length and K is the beam size with the derivation in Section 6.

4 Experimental Setup

We introduce the MT datasets and pre-processing settings Section 4.1. In Section 4.2, we provide details around the model hyper-parameters, training and inference settings. In Section 4.3, we present our evaluation metrics.

4.1 Data and Pre-processing

Datasets We conducted MT experiments with datasets listed in Table 1, including WMT'22 Livonian–English (Liv–En), Asian Language Treebank (ALT), IWSLT'15 Vietnamese–English (Vi–En), WMT'16 Romanian–English (Ro–En), WMT'15 Finnish–English (Fi–En), and WMT'14 German–English (De–En) datasets. ALT is a multi-way parallel dataset containing data in English and other Asian languages including Filipino (Fil), Indonesian (Id), Japanese (Ja), Malay

Algorithm 2: InnerBeamSearch

Data: Beam width K , max length T , $toks$
Result: Next word list

```
1 Initialization:
2  $B_0 \leftarrow \{(0, toks)\}$ ;
3 for  $t \leftarrow 1$  to  $T$  do
4    $B_t \leftarrow \emptyset$ ;
5   foreach  $s \in B_{t-1}$  do
6     if  $s$  reaches  $\_$  or  $\langle eos \rangle$  then
7        $B_t.append(s)$ ;
8       continue;
9     foreach  $s' \in Decoder(s[1])$  do
10       $score, toks = s'$ ;
11       $score \leftarrow s[0] + score$ ;
12       $toks \leftarrow s[1] + toks$ ;
13       $B_t.append((score, toks))$ ;
14   Sort  $B_t$  by scores in descending order.;
15    $B_t \leftarrow B_t[: K]$ 
16  $W = \{\}$ ;
17 foreach  $s \in B_T$  do
18    $score, toks = s$ ;
19    $word = detokenize(toks)$ ;
20   if  $word \notin W$  then
21      $W[word] = (score, toks)$ 
22   else
23      $W[word][0] += score$ 
24 return  $list(W.items())$ 
```

Dataset	Train	Valid	Test
WMT'22 Liv-En	1, 127	586	856
ALT Asian Langs-En	18k	1, 000	1, 018
IWSLT'15 Vi-En	133k	1, 553	1, 268
WMT'16 Ro-En	612k	1, 999	1, 999
WMT'15 Fi-En	1.8M	1, 500	1, 370
WMT'14 De-En	4.5M	45, 781	3, 003

Table 1: Statistics of the datasets.

(Ms), Vietnamese (Vi), and simplified Chinese (Zh). We used the ALT-standard-split tool² to split the dataset into train, validation, and test sets.

Data Pre-processing We performed word tokenization on all data. We applied Juman++ (Tolmachev et al., 2018) to data in Japanese, Stanford-tokenizer (Manning et al., 2014) to data in Chinese, and Moses tokenizer (Koehn et al., 2007) to data in other languages. We normalized Romanian data and removed diacritics following previous work (Sennrich et al., 2016a). We prepared the WMT'14 English-German dataset using a data cleaning and normalization tool from Fairseq.³

²www2.nict.go.jp/astrec-att/member/mutiyama/ALT

³github.com/facebookresearch/fairseq/blob/main/examples/translation/

We applied subword tokenization to each translation direction separately. For source or target language, we trained a subword tokenizer with a subword vocabulary of $8k$ on the monolingual corpus from the training set. The vocabulary size is computed by the VOLT algorithm (Xu et al., 2021). For languages in WMT'22, ALT and IWSLT'15, they are $7k$ to $8k$, and for the remaining datasets they are $10k$ to $11k$. We used $8k$ for consistency. We applied a widely adopted toolkit⁴ to train BPE-dropout tokenizers with a dropout rate of 0.2 for the generation of regularized data and train BPE tokenizers for the generation of non-regularized data. The dropout rate is selected through hyperparameter grid search from 0.1 to 0.5 with steps of 0.1, where we found 0.2 usually optimal and rate ≥ 0.3 resulted in unstable training.

4.2 NMT Settings

Model We used the Fairseq framework (Ott et al., 2019). We refer model settings in previous works (Rubino et al., 2020; Provilkov et al., 2020). For WMT'22, ALT and IWSLT'15 datasets, we used 1 attention head, 6 decoder layers, and 4 or 6 encoder layers (4 layers only for En \leftarrow \rightarrow Fil and Ja \rightarrow En) and FFN dim of 512. For other datasets we used the standard transformer base architecture (Vaswani et al., 2017). We set dropout and attention dropout rates to 0.1. We applied layer normalization (Lei Ba et al., 2016; Mao et al., 2023) for both the encoder and decoder.

Training We set the batch size to 3,072 tokens for sentence in the source language and used eight GPUs, resulting in 25k tokens per batch. We used the Adam optimizer (Kingma and Ba, 2014) with $\beta_1 = 0.9$, $\beta_2 = 0.98$, and $\epsilon = 10^{-9}$. We used warmup and linear decay for the learning rate (Vaswani et al., 2017), with 4k warm-up steps, an initial learning rate of $1.7 * 10^{-7}$ and a final learning rate of $5 * 10^{-4}$. We used label smoothing for the cross entropy loss with $\epsilon_{ls} = 0.1$ (Szegedy et al., 2015). We calculated the loss on the validation set after each epoch and applied early stopping when no improvement was observed for 10 epochs.

SubMerge led to better word-level perplexities than traditional beam search and higher BLEU and chrF++ scores, often achieving statistically significant improvements.

⁴[prepare-wmt14en2de.sh](https://github.com/google/sentencepiece)

⁴github.com/google/sentencepiece

	Word Perplexity ↓		BLEU ↑		chrF++ ↑	
	Beam Search	SubMerge	Beam Search	SubMerge	Beam Search	SubMerge
Low-Resource Scenario						
WMT’22 Liv→En	5.93	3.43	1.52	2.04 ^{+0.5}	18.85	19.45 ^{+0.6}
WMT’22 En→Liv	19.39	6.88	2.70	3.21 ^{+0.5}	19.14	19.41 ^{+0.3}
ALT Fil→En	12.68	4.59	31.10	31.82 ^{+0.7}	57.98	59.17 ^{+1.2}
ALT En→Fil	9.56	4.14	30.20	31.14 ^{+0.9}	59.64	60.14 ^{+0.5}
ALT Id→En	17.91	5.91	27.35	28.73 ^{+1.4}	53.61	56.39 ^{+2.8}
ALT En→Id	16.44	4.91	33.63	34.19 ^{+0.6}	63.14	63.89 ^{+0.8}
ALT Ja→En	24.90	7.79	15.07	15.26 ^{+0.2}	45.07	45.46 ^{+0.4}
ALT En→Ja	6.55	3.69	14.38	14.59 ^{+0.2}	27.92	29.02 ^{+1.1}
ALT Ms→En	11.28	4.33	31.86	32.16 ^{+0.3}	59.01	60.09 ^{+1.1}
ALT En→Ms	12.82	4.18	38.83	39.28 ^{+0.5}	66.25	66.91 ^{+0.7}
ALT Vi→En	17.21	6.14	23.64	24.97 ^{+1.3}	52.32	52.93 ^{+0.6}
ALT En→Vi	8.64	3.52	27.35	27.64 ^{+0.3}	53.66	53.82 ^{+0.2}
ALT Zh→En	23.11	7.81	13.92	14.31 ^{+0.4}	43.54	44.43 ^{+0.9}
ALT En→Zh	13.61	6.76	9.03	9.87 ^{+0.8}	22.76	23.25 ^{+0.5}
Middle- and High-Resource Scenario						
IWSLT’15 Vi→En	14.41	5.62	27.87	28.43 ^{+0.6}	48.62	50.59 ^{+2.0}
IWSLT’15 En→Vi	7.98	3.39	28.08	28.16 ^{+0.1}	49.27	50.18 ^{+0.9}
WMT’16 Ro→En	7.44	3.22	33.85	33.77 ^{-0.1}	58.75	59.07 ^{+0.3}
WMT’16 En→Ro	6.78	3.11	34.35	34.50 ^{+0.1}	58.66	58.89 ^{+0.2}
WMT’15 Fi→En	11.27	4.27	18.95	18.88 ^{-0.1}	47.24	47.55 ^{+0.3}
WMT’15 En→Fi	22.52	7.81	16.51	16.65 ^{+0.1}	47.66	47.97 ^{+0.3}
WMT’14 De→En	10.33	3.90	28.85	28.94 ^{+0.1}	55.99	56.52 ^{+0.5}
WMT’14 En→De	12.74	4.64	24.69	24.83 ^{+0.1}	52.68	52.77 ^{+0.1}

Table 2: Results of Subword Regularized Models. Statistical significance $p < 0.01$ is indicated by * against Beam Search. SubMerge consistently improves over the Beam Search baseline in most directions. Word perplexity results represent the ability to accurately estimate sentence probability rather than fluency.

Inference We selected the checkpoint with the best loss on the validation set. We used beam search and SubMerge with a beam size of 4 without additional normalization techniques, such as length penalty or temperature sampling (Dong et al., 2022).

4.3 Evaluation Metrics

We report word perplexity on generated translations to compare the probabilities assigned to generations by models. To evaluate translation quality, we report BLEU using sacreBLEU (Post, 2018),⁵ chrF++ (Popović, 2017),⁶ and BLEURT (Appendix A). We performed paired bootstrap resampling for statistical significance tests (Koehn, 2004).

The word perplexity is calculated as follows. We first evaluate the negative log probability of the generated sentences for models using SubMerge by:

$$s_{score} = - \sum_i \log P_{\theta}(word_i), \quad (7)$$

⁵BLEU+c.mixed+l.en-lang+#.l+s.exp+tok.l3a+v.l.5.1

⁶github.com/m-popovic/chrF with c6w2F0.4. Similar trends were observed using different chrF settings.

and models with beam search by:

$$s_{score} = - \sum_i \log P_{\theta}(tok_i). \quad (8)$$

We evaluated the average word perplexity by

$$w_{ppl} = \exp\left(\frac{1}{N} s_{score}\right), \quad (9)$$

where N is the number of words. We evaluated the word perplexity based on the generated hypothesis rather than the reference. This reflects the actual scenario in generation tasks where we dynamically generate the next token (word) conditioned on what the model has generated instead of on the ground truth. Nevertheless, word perplexity is a conditional probability dependent on not only the input but also the parameters in the model. Therefore, the perplexity results must always be considered along with model-independent metrics such as BLEU scores.

5 Translation Quality Results

The results for subword regularized models are shown in Table 2.

Word Perplexity We observed that word perplexity results improved substantially in the regularized models in contrast to the tiny gap (0.5%) reported in the non-regularized models (Chirkova et al., 2023) and in our analysis shown in Section 7.5. This is due to the fact that multiple tokenizations for one word appeared during training, which acts as a label-smoothing function on multiple correct next tokens. Therefore, the probability weight is distributed across multiple subwords thus, it becomes necessary to incorporate the marginal likelihood. It is worth noting that here word perplexity represents the precision of probability estimation rather than fluency or quality of the output.

Translation Quality We also found translation quality improved, especially in low-resource scenarios where the average BLEU score improvement is 0.6, whereas in higher resource scenarios, it is 0.3. We also observed consistent improvement in the chrF++ score. While only one translation direction among higher resource directions is statistically significant, 8 out of 12 low-resource directions see statistically significant improvements. Furthermore, we observed that the improvement is greater for languages where words contain more subwords on average ($T_{subword}$). In the ALT dataset, each Japanese word contains an average of 1.59 subwords, resulting in a modest improvement of only 0.2 BLEU. In contrast, Filipino has a $T_{subword}$ of 2.16, leading to an improvement of 0.9 BLEU.

6 Efficiency

We show the theoretical analysis of time complexity as well as running time results of SubMerge comparing with beam search.

Time Complexity Let K denote beam size and T_{word} denote the number of words in the sentence. In the outer beam search Algorithm 1, the loop in line 3 contains at most T steps, and line 5 contains at most K steps. Therefore, the time complexity of the SubMerge is $O(T_{word} * K * O(InnerBeamSearch()))$.

In the inner beam search Algorithm 2, line 3 contains at most $T_{subword}$ steps, which is the number of subwords within the word boundary. Line 5 contains at most K steps, and line 9 contains at most K steps because each beam yields maximum K candidates by selecting top- K probable tokens. Therefore, the time complexity of Algorithm 2 is $O(T_{subword} * K * K)$.

The overall time complexity of SubMerge is $O(\sum_i(T_{subword}) * K^3) = O(T * K^3)$ which is K times slower than that of beam search which is $O(T * K^2)$.

Inference Time We compared the running times in the IWSLT’15 En→Vi direction using $K = 4$ extracted from the log data reported by the Fairseq framework. SubMerge took 1,665 seconds to generate 1,268 sentences, whereas beam search took 303 seconds, showing SubMerge is approximately 5.5 times slower. We set the batch size to 1 because the current SubMerge implementation does not yet support batch processing.

7 Analysis

We investigate the effect of different beam sizes on the algorithm in Section 7.1. Section 7.2 explores using a sampling algorithm as the inner search algorithm. Section 7.3 and Section 7.4 respectively analyze the impact of the training set size and the dropout rate. Section 7.5 show conditions in which SubMerge is effective.

7.1 Assessing Beam Sizes Variants

Figures 3 and 4 show the word perplexities and BLEU scores of using different beam sizes for both non-regularized models and subword regularized models, comparing beam search and SubMerge.

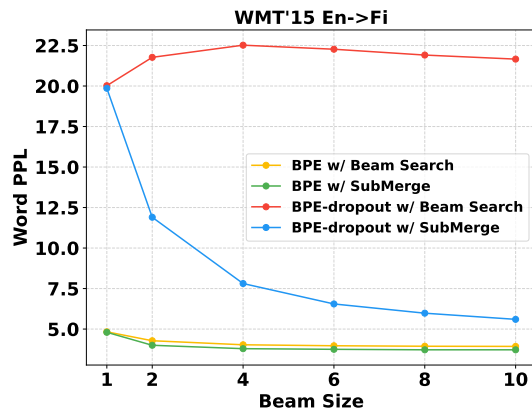


Figure 3: Word perplexity results using different beam sizes on the WMT’15 En→Fi direction.

We observed that as we increased the beam size, the word perplexity dropped sharply for BPE-dropout with SubMerge. When using a large beam size such as 10, it achieved comparable results to non-regularized models trained on one-best tokenization. Nevertheless, SubMerge does not yet accumulate as large a proportion of the

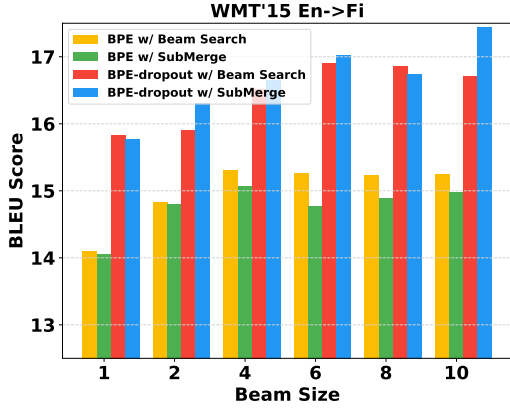


Figure 4: BLEU results using different beam sizes on the WMT'15 En→Fi direction.

probability distribution as using a non-regularized model. Since the training is on multiple segmentations, it certainly comes closer than when using beam search. For non-regularized models, combining equivalent paths for perplexity estimation also proved to be effective. We also observed that increasing beam size can lead to translation quality improvement for the SubMerge method. However, this is not the case for all directions (Cohen and Beck, 2019) and we put the full results using different beam sizes for all datasets in Appendix C.

7.2 Inner Search Algorithm Variants

We replaced the inner beam search with the sampling algorithm as shown in Algorithm 3. In the algorithm, Q is the queue that contains possible subword tokenizations of the next word. The sampling algorithm selects the next token tok_j in line 10 for each ongoing sample s according to the probability distribution of subwords in the target vocabulary outputted by the softmax function after the decoder. The current s is updated for both the score and the string. We call this pure sampling because we did not add sampling temperature, top- k or top- p filtering. We perform the merging post-processing the same as the inner beam search.

The word perplexity results are shown in Figure 5. For the sampling algorithm, we sampled n^2 tokenizations (where n is the beam size) in the inner loop and for each path, we started with the same historical information and selected the next subword according to the probability distribution until we reached the beginning of the next word. We then perform the same merging post-processing. However, we observed that the perplexity was higher than n -best tokenizations. This

is because the sampling process could easily get lost at some step by selecting a token in the long tail with a very low probability.

Algorithm 3: InnerSampling

Data: Sample times K , max length T , $toks$
Result: Next word list

```

1 Initialization:
2  $s_0 \leftarrow \{(0, toks)\}$ ;
3  $Q \leftarrow \emptyset$ ;
4 for  $i \leftarrow 1$  to  $K$  do
5    $s \leftarrow s_0$ ;
6   for  $j \leftarrow 1$  to  $T$  do
7     if  $s$  reaches  $\_or < eos >$  then
8        $Q.append(s)$ ;
9       break;
10    Sample  $tok_j$  from  $Decoder(s[1])$ ;
11    Update  $s$  using  $tok_j$ ;
12 Sort  $Q$  by scores in descending order.;
13  $W = \{\}$ ;
14 foreach  $s \in Q$  do
15    $score, toks = s$ ;
16    $word = detokenize(toks)$ ;
17   if  $word \notin W$  then
18      $W[word] = (score, toks)$ 
19   else
20      $W[word][0] += score$ 
21 return  $list(W.items())$ 

```

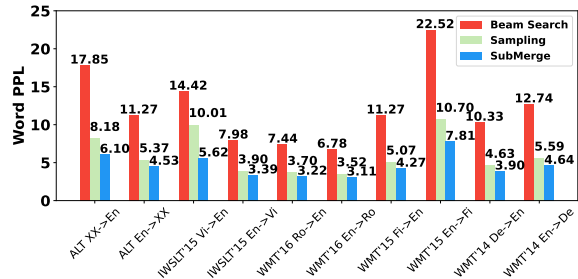


Figure 5: Word perplexity results comparing BPE-dropout with beam search to two variants of SubMerge: using either sampling as the inner search function or beam search.

7.3 Assessing Training Set Sizes

SubMerge is effective in extremely low-resource scenarios, as shown in Figure 6. We reported BLEU scores using beam search and SubMerge during decoding for models trained on 1k to 18k parallel sentences. SubMerge consistently outperformed beam search across training set sizes. Moreover, the BLEU improvement reached approximately 3.4 using only 1k data. This observation reveals the potential of SubMerge to be used in domain adaptation scenarios with limited data.

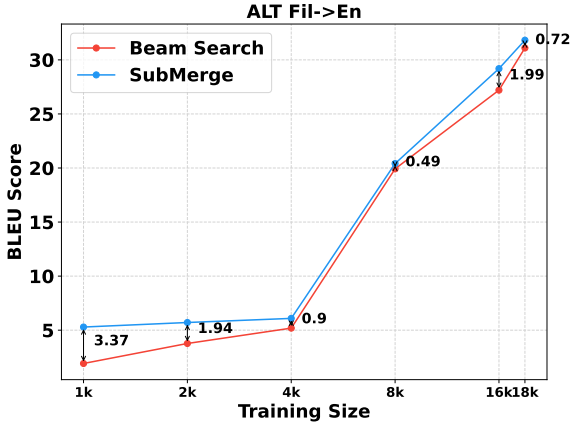


Figure 6: Translation quality using different sizes of training data. The x-axis is logarithmized.

7.4 Impact of Dropout Rates

Using a lower dropout rate in BPE-dropout yielded lower word perplexity and higher BLEU scores in higher resource scenarios, as shown in Table 3. When the dropout rate is low, the randomness of subword segmentation for a given word also decreases, leading to reduced variability in the training data and, concurrently, a diminished range of choices during the inference process. In the context of low-resource scenarios, reduced variability implies diminished data augmentation, which can adversely affect the model’s generalization capability. Conversely, in higher resource settings, decreased variability signifies reduced noise, potentially enhancing model performance.

Dropout Rate	Word PPL ↓		BLEU ↑	
	0.1	0.2	0.1	0.2
ALT Others→En	4.69	6.10	22.06	24.54
ALT En→Others	4.16	4.53	24.75	26.12
IWSLT’ 15 Vi→En	3.09	5.62	30.03	28.43
IWSLT’ 15 En→Vi	2.56	3.39	29.61	28.16
WMT’ 16 Ro→En	2.34	3.22	34.75	33.77
WMT’ 16 En→Ro	2.21	3.11	35.39	34.50
WMT’ 15 Fi→En	3.25	4.27	18.87	18.88
WMT’ 15 En→Fi	4.94	7.81	16.64	16.65
WMT’ 14 De→En	2.86	3.90	29.70	28.94
WMT’ 14 En→De	3.15	4.64	24.94	24.83

Table 3: Results of SubMerge for models trained on BPE-dropout data with different dropout rates.

7.5 Does SubMerge Work on Non-regularized Models?

In short, No. We explored whether the proposed SubMerge method is applicable to non-regularized models using deterministic BPE tokenization. Ta-

	Word PPL ↓		BLEU ↑	
	BeamSearch	SubMerge	BeamSearch	SubMerge
WMT’ 22 Liv→En	3.60	3.37	0.36	0.44
WMT’ 22 En→Liv	5.22	4.55	0.64	0.90
ALT Others→En	6.02	5.60	15.73	15.40
ALT En→Others	4.90	4.77	18.06	17.82
IWSLT’ 15 Vi→En	2.95	2.79	24.34	25.63
IWSLT’ 15 En→Vi	2.43	2.42	25.09	24.86
WMT’ 16 Ro→En	2.14	2.11	32.05	31.70
WMT’ 16 En→Ro	2.00	1.98	32.98	32.85
WMT’ 15 Fi→En	2.85	2.76	17.08	16.94
WMT’ 15 En→Fi	4.03	3.79	15.30	15.06
WMT’ 14 De→En	2.39	2.40	30.18	30.04
WMT’ 14 En→De	2.45	2.36	25.88	25.71

Table 4: Results of non-regularized models trained on data using BPE tokenizer. We show the averaged results in En→XX and XX→En directions for the ALT dataset.

ble 4 presents word perplexities and BLEU scores on non-regularized models using beam search or SubMerge as the decoding algorithm.

We observed lower word perplexity using SubMerge compared to using beam search. However, the improvement is not as significant (approximately 6%) as the improvement achieved by SubMerge for subword regularized models. This is consistent with our expectations. Models were trained on a single tokenization for each training word, so one tokenization accumulates the most probability weight. For the non-regularized model, results show the translation quality of SubMerge is not as good as that of beam search. Therefore, the proposed SubMerge method is only applicable to subword regularized models in the NMT task.

For other tasks, such as question answering, the word perplexity is greater because the task is less structured than MT, where the source sentence is a highly limiting constraint. For less constrained tasks, it is possible that SubMerge will improve the performance of even non-regularized models. We leave this for future work to explore.

8 Related Work

SubMerge is designed for decoding with text generation models for which likely tokenization probabilities diverge drastically from sentence probabilities. In other words, there are multiple tokenizations for one target sentence, and the probability distribution is splintered among them. Our objective is to enhance the inference algorithm on the target side. On the source side, merging probabilities of multiple tokenizations for a single source sentence has been shown to improve translation performance in low-resource scenar-

ios (Takase et al., 2022). Although we only experimented on models trained on data segmented by BPE-dropout (Provilkov et al., 2020), it also works for SentencePiece Regularization (Kudo, 2018a), MaxMatch-Dropout (Hiraoka, 2022) and NMT models with multiple subword segmenters (Kambhatla et al., 2022). On the other hand, NMT models trained on sentences segmented by deterministic segmenter only benefit from marginal likelihood estimation in out-of-domain data or long words (Cao and Rimell, 2021; Chirkova et al., 2023). Deterministic subword segmentation includes not only subword-level methods such as WordPiece (Schuster and Nakajima, 2012), BPE (Sennrich et al., 2016b), SentencePiece (Kudo and Richardson, 2018), dynamic programming encoding (He et al., 2020), BERT-Seg (Song et al., 2022), but also byte-level (Shaham and Levy, 2021), character-level (Tay et al., 2021), word-level (Mikolov et al., 2013), and hybrid word-character methods (Luong and Manning, 2016).

Marginal likelihood estimation can be implemented in two ways: sampling and dynamic programming. Sampling methods include summing over n -best tokenizations (Cao and Rimell, 2021) or important tokenizations (Chirkova et al., 2023). Sampling can be easily applied to any generation model. However, a manageable number of tokenizations cannot precisely estimate the probability of sentences with an exponentially large number of tokenizations, which is the case during the inference of the subword regularized models. On the other hand, dynamic programming can handle an exponentially large number of tokenizations by merging the same historical states, as introduced in *sequence modeling via segmentations* (Wang et al., 2017) and applied in the mixed-character-subword models (He et al., 2020; Meyer and Buys, 2023). However, they merge the historical states by approximating the previous output by character-level data. That is, after the decoder generates one subword, it is split into characters and fed to the decoder. This is not applicable to pure subword models. Based on the property that each word is individually segmented in BPE-dropout (Provilkov et al., 2020), we obtain n -best tokenizations within a small search space and treat the best tokenization of each word the historical state, taking advantage of both marginal likelihood estimation methods.

9 Conclusion and Future Work

We propose SubMerge to estimate the marginal likelihood of the next word by merging equivalent subword tokenizations during the inference of subword regularized models. Results demonstrate a significant improvement in word perplexity estimation and translation quality improvement in terms of BLEU and chrF++ scores, especially in low-resource scenarios.

Current inference algorithms are mostly based on conditional probability, which is a short-term value function. For future work of inference, we suggest aligning the value function towards evaluation metrics and human preference through reinforcement learning, where models are more aware of longer-term rewards.

Limitations

We did not experiment with common techniques in the beam search and SubMerge, such as length penalty. This is because we use a nested beam search, and the way to define the length (whether to use the number of tokens or the number of words) may differ from the definition in a traditional beam search. However, combining SubMerge with such techniques could be valuable for further work.

The word perplexity results reported in this paper are on the generated texts rather than reference texts. They do not correlate with fluency or translation quality, and we only use them to report how much of the probability weight of a model is being used during decoding, which is still useful.

We use the SentencePiece tool for the current implementation of BPE and BPE-dropout algorithms. Therefore, the SubMerge implementation is also based on the format of this specific tool, which uses “_” (U+2581) to represent the beginning of a new word. However, other tools may use “@@” at the end of a subword to indicate that the current word has not ended yet. Therefore, the implementation of SubMerge may be slightly different in terms of ending conditions in the inner beam search.

We did not experiment on large-scale datasets (e.g., datasets with more than 100M parallel sentences). Reasons include 1) computational budget limitations and 2) the goal is verifying the algorithm rather than developing systems. We assume that the improvement will be marginal in high-resource scenarios.

References

- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv e-prints*, page arXiv:1409.0473, September.
- Cao, Kris and Laura Rimell. 2021. You should evaluate your language model on marginal likelihood over tokenisations. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2104–2114, Online and Punta Cana, Dominican Republic, November. Association for Computational Linguistics.
- Cherry, Colin, George Foster, Ankur Bapna, Orhan Firat, and Wolfgang Macherey. 2018. Revisiting character-based neural machine translation with capacity and compression. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4295–4305, Brussels, Belgium, October–November. Association for Computational Linguistics.
- Chirkova, Nadezhda, Germán Kruszewski, Jos Rozen, and Marc Dymetman. 2023. Should you marginalize over possible tokenizations? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1–12, Toronto, Canada, July. Association for Computational Linguistics.
- Cohen, Eldan and Christopher Beck. 2019. Empirical analysis of beam search performance degradation in neural sequence models. In Chaudhuri, Kamalika and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1290–1299. PMLR, 09–15 Jun.
- Costa-jussà, Marta R. and José A. R. Fonollosa. 2016. Character-based neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 357–361, Berlin, Germany, August. Association for Computational Linguistics.
- Dong, Chenhe, Yinghui Li, Haifan Gong, Miaoxin Chen, Junxin Li, Ying Shen, and Min Yang. 2022. A survey of natural language generation. *ACM Comput. Surv.*, 55(8), dec.
- Gehring, Jonas, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. 2017. Convolutional sequence to sequence learning. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML’17, page 1243–1252. JMLR.org.
- Gupta, Rohit, Laurent Besacier, Marc Dymetman, and Matthias Gallé. 2019. Character-based nmt with transformer.
- He, Xuanli, Gholamreza Haffari, and Mohammad Norouzi. 2020. Dynamic programming encoding for subword segmentation in neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3042–3051, Online, July. Association for Computational Linguistics.
- Hiraoka, Tatsuya. 2022. MaxMatch-dropout: Subword regularization for WordPiece. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4864–4872, Gyeongju, Republic of Korea, October. International Committee on Computational Linguistics.
- Kalchbrenner, Nal and Phil Blunsom. 2013. Recurrent continuous translation models. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1700–1709, Seattle, Washington, USA, October. Association for Computational Linguistics.
- Kambhatla, Nishant, Logan Born, and Anoop Sarkar. 2022. Auxiliary subword segmentations as related languages for low resource multilingual translation. In *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, pages 131–140, Ghent, Belgium, June. European Association for Machine Translation.
- Kim, Yoon, Yacine Jernite, David Sontag, and Alexander Rush. 2016. Character-aware neural language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 30(1), Mar.
- Kingma, Diederik P. and Jimmy Ba. 2014. Adam: A method for stochastic optimization.
- Koehn, Philipp, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Prague, Czech Republic, June. Association for Computational Linguistics.
- Koehn, Philipp. 2004. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain, July. Association for Computational Linguistics.
- Kudo, Taku and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium, November. Association for Computational Linguistics.

- Kudo, Taku. 2018a. Subword regularization: Improving neural network translation models with multiple subword candidates. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia, July. Association for Computational Linguistics.
- Kudo, Taku. 2018b. Subword regularization: Improving neural network translation models with multiple subword candidates.
- Lei Ba, Jimmy, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. Layer Normalization. *arXiv e-prints*, page arXiv:1607.06450, July.
- Ling, Wang, Isabel Trancoso, Chris Dyer, and Alan W Black. 2015. Character-based neural machine translation.
- Luong, Minh-Thang and Christopher D. Manning. 2016. Achieving open vocabulary neural machine translation with hybrid word-character models.
- Luong, Thang, Ilya Sutskever, Quoc Le, Oriol Vinyals, and Wojciech Zaremba. 2015. Addressing the rare word problem in neural machine translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 11–19, Beijing, China, July. Association for Computational Linguistics.
- Manning, Christopher, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. The Stanford CoreNLP natural language processing toolkit. In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60, Baltimore, Maryland, June. Association for Computational Linguistics.
- Mao, Zhuoyuan, Raj Dabre, Qianying Liu, Haiyue Song, Chenhui Chu, and Sadao Kurohashi. 2023. Exploring the impact of layer normalization for zero-shot neural machine translation. In Rogers, Anna, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1300–1316, Toronto, Canada, July. Association for Computational Linguistics.
- Meyer, Francois and Jan Buys. 2023. Subword segmental machine translation: Unifying segmentation and target sentence generation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 2795–2809, Toronto, Canada, July. Association for Computational Linguistics.
- Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space.
- Ott, Myle, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Popović, Maja. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels, October. Association for Computational Linguistics.
- Provlkov, Ivan, Dmitrii Emelianenko, and Elena Voita. 2020. BPE-dropout: Simple and effective subword regularization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1882–1892, Online, July. Association for Computational Linguistics.
- Pu, Amy, Hyung Won Chung, Ankur P Parikh, Sebastian Gehrmann, and Thibault Sellam. 2021. Learning compact metrics for mt. In *Proceedings of EMNLP*.
- Rubino, Raphael, Benjamin Marie, Raj Dabre, Atushi Fujita, Masao Utiyama, and Eiichiro Sumita. 2020. Extremely low-resource neural machine translation for asian languages. *Machine Translation*, 34(4):347–382.
- Schuster, M. and K. Nakajima. 2012. Japanese and korean voice search. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5149–5152.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2016a. Edinburgh neural machine translation systems for WMT 16. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 371–376, Berlin, Germany, August. Association for Computational Linguistics.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2016b. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August. Association for Computational Linguistics.
- Shaham, Uri and Omer Levy. 2021. Neural machine translation without embeddings. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics*:

- Human Language Technologies*, pages 181–186, Online, June. Association for Computational Linguistics.
- Song, Haiyue, Raj Dabre, Zhuoyuan Mao, Chenhui Chu, and Sadao Kurohashi. 2022. BERTSeg: BERT based unsupervised subword segmentation for neural machine translation. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 85–94, Online only, November. Association for Computational Linguistics.
- Sutskever, Ilya, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In Ghahramani, Z., M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 3104–3112. Curran Associates, Inc.
- Szegedy, Christian, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2015. Re-thinking the inception architecture for computer vision.
- Takase, Sho, Tatsuya Hiraoka, and Naoaki Okazaki. 2022. Single model ensemble for subword regularized models in low-resource machine translation. In Muresan, Smaranda, Preslav Nakov, and Aline Villavicencio, editors, *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2536–2541, Dublin, Ireland, May. Association for Computational Linguistics.
- Tay, Yi, Vinh Q. Tran, Sebastian Ruder, Jai Gupta, Hyung Won Chung, Dara Bahri, Zhen Qin, Simon Baumgartner, Cong Yu, and Donald Metzler. 2021. Charformer: Fast character transformers via gradient-based subword tokenization.
- Tolmachev, Arseny, Daisuke Kawahara, and Sadao Kurohashi. 2018. Juman++: A morphological analysis toolkit for scriptio continua. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 54–59, Brussels, Belgium, November. Association for Computational Linguistics.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Guyon, I., U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Wang, Chong, Yining Wang, Po-Sen Huang, Abdelrahman Mohamed, Dengyong Zhou, and Li Deng. 2017. Sequence modeling via segmentations.
- Xu, Jingjing, Hao Zhou, Chun Gan, Zaixiang Zheng, and Lei Li. 2021. Vocabulary learning via optimal transport for neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7361–7373, Online, August. Association for Computational Linguistics.

A BLEURT Results

Table 5 shows BLEURT score results using the BLEURT-20 model (Pu et al., 2021). We can observe a similar trend with other metrics such as BLEU or chrF++, where the improvement is large in low-resource directions and comparable in higher-resource directions.

B Comparing with Non-Subword Models

We trained character-based and word-based models on IWSLT’15 Vi–En and WMT’16 Ro–En datasets and showed inferior performance compared to subword-based models using SubMerge as shown in Table 6. This conclusion is aligned with that in previous paper (Kudo, 2018b) and report.⁷

C Full Results for Different Beam Sizes

Tables 7, 8 and 9 show negative sentence log probability, word perplexity and BLEU scores for different beam sizes. The conclusions still remain the same where SubMerge improved probability estimation precision, which however did not bring translation quality improvement.

	BLEURT ↑	
	Beam Search	SubMerge
<i>Low-Resource Scenario</i>		
WMT’22 Liv→En	17.40	17.47
WMT’22 En→Liv	42.74	42.40
ALT Fil→En	55.35	56.70
ALT En→Fil	47.95	47.78
ALT Id→En	51.65	53.91
ALT En→Id	56.72	57.10
ALT Ja→En	41.54	41.88
ALT En→Ja	27.02	26.86
ALT Ms→En	56.87	57.83
ALT En→Ms	59.32	59.44
ALT Vi→En	49.61	50.97
ALT En→Vi	44.79	45.11
ALT Zh→En	40.95	41.38
ALT En→Zh	28.19	29.19
<i>Middle- and High-Resource Scenario</i>		
IWSLT’15 Vi→En	51.75	52.57
IWSLT’15 En→Vi	47.13	47.46
WMT’16 Ro→En	61.52	61.35
WMT’16 En→Ro	52.08	51.76
WMT’15 Fi→En	57.12	56.84
WMT’15 En→Fi	53.67	53.41
WMT’14 De→En	60.01	59.82
WMT’14 En→De	54.97	54.27

Table 5: BLEURT Results of Subword Regularized Models.

⁷github.com/google/sentencepiece/blob/master/doc/experiments.md

Models	Vi→En	En→Vi	Ro→En	En→Ro
Char-based	24.72	27.04	30.45	29.82
Word-based	21.40	25.24	26.33	25.67
Subword-based	28.43	28.16	33.77	34.50

Table 6: BLEU score results comparing different models on IWSLT’15 Vi–En and WMT’16 Ro–En datasets.

	IWSLT’15 Vi→En	IWSLT’15 En→Vi	WMT’16 Ro→En	WMT’16 En→Ro	WMT’15 Fi→En	WMT’15 En→Fi	WMT’14 De→En	WMT’14 En→De
<i>BeamSize=1</i>								
BPE w/ Beam Search	25.67	24.37	19.86	19.08	23.64	23.30	21.91	21.01
BPE w/ SubMerge	24.82	24.06	22.87	20.62	23.40	23.07	28.98	20.25
BPE-dropout w/ Beam Search	30.51	31.59	28.11	27.68	32.35	31.22	31.77	33.45
BPE-dropout w/ SubMerge	30.16	31.14	27.83	27.56	31.81	30.99	-	32.57
<i>BeamSize=2</i>								
BPE w/ Beam Search	23.54	22.52	18.75	17.99	21.83	21.21	19.90	20.11
BPE w/ SubMerge	22.15	22.18	18.42	18.96	21.36	20.69	18.71	19.32
BPE-dropout w/ Beam Search	29.74	30.58	27.93	27.22	31.84	30.55	31.25	33.03
BPE-dropout w/ SubMerge	25.83	27.15	23.95	23.45	27.52	26.75	26.59	28.72
<i>BeamSize=3</i>								
BPE w/ Beam Search	22.95	21.90	18.39	17.53	21.18	20.43	19.31	19.75
BPE w/ SubMerge	21.33	21.66	18.12	18.17	20.66	19.95	18.44	18.97
BPE-dropout w/ Beam Search	29.56	30.28	27.87	26.93	31.63	30.00	31.05	32.74
BPE-dropout w/ SubMerge	23.59	24.70	21.56	21.16	30.45	24.69	24.21	26.41
<i>BeamSize=4</i>								
BPE w/ Beam Search	22.59	21.61	18.21	17.31	20.80	20.10	19.08	19.54
BPE w/ SubMerge	20.97	21.31	17.84	17.20	20.39	19.75	20.16	18.82
BPE-dropout w/ Beam Search	29.52	30.06	27.79	26.79	31.46	29.71	30.88	32.57
BPE-dropout w/ SubMerge	22.57	23.65	20.35	19.86	23.82	23.72	22.91	25.08
<i>BeamSize=5</i>								
BPE w/ Beam Search	22.43	21.42	18.07	17.18	20.50	19.81	18.88	19.43
BPE w/ SubMerge	20.66	21.14	17.70	17.64	20.16	19.53	19.33	18.70
BPE-dropout w/ Beam Search	29.42	29.82	27.67	26.71	31.56	29.52	30.72	32.41
BPE-dropout w/ SubMerge	22.38	22.74	19.39	19.05	25.77	22.76	22.02	24.11
<i>BeamSize=6</i>								
BPE w/ Beam Search	22.21	21.31	18.01	17.06	20.42	19.67	18.77	19.37
BPE w/ SubMerge	20.46	20.96	17.69	16.97	20.09	19.41	19.18	18.62
BPE-dropout w/ Beam Search	29.42	29.65	27.67	26.65	31.20	29.41	30.66	32.36
BPE-dropout w/ SubMerge	21.75	22.31	18.82	18.38	22.61	22.23	21.41	23.49
<i>BeamSize=8</i>								
BPE w/ Beam Search	21.95	21.13	17.85	16.87	20.21	19.45	18.66	19.25
BPE w/ SubMerge	20.26	20.77	17.53	17.44	19.80	19.18	18.74	18.51
BPE-dropout w/ Beam Search	29.08	29.51	27.57	26.57	31.19	29.18	30.49	32.14
BPE-dropout w/ SubMerge	21.20	22.07	18.28	17.87	21.80	21.65	20.83	22.80
<i>BeamSize=10</i>								
BPE w/ Beam Search	21.77	20.93	17.75	16.72	20.06	19.24	18.60	19.16
BPE w/ SubMerge	20.08	20.62	17.41	17.06	19.69	19.02	18.41	18.39
BPE-dropout w/ Beam Search	28.84	29.39	27.46	26.47	30.69	29.02	30.40	32.03
BPE-dropout w/ SubMerge	20.82	21.56	17.84	17.44	22.03	21.12	20.37	22.26

Table 7: Negative sentence log probability of the generated hypothesis using different beam sizes.

	IWSLT'15 Vi→En	IWSLT'15 En→Vi	WMT'16 Ro→En	WMT'16 En→Ro	WMT'15 Fi→En	WMT'15 En→Fi	WMT'14 De→En	WMT'14 En→De
<i>BeamSize=1</i>								
BPE w/ Beam Seach	3.34	2.67	2.27	2.12	3.19	4.83	2.64	2.58
BPE w/ SubMerge	3.30	2.67	2.55	2.26	3.20	4.80	3.20	2.52
BPE-dropout w/ Beam Seach	4.42	3.42	3.20	2.99	4.79	8.32	4.19	4.59
BPE-dropout w/ SubMerge	4.42	3.41	3.18	3.00	4.76	8.28	-	4.47
<i>BeamSize=2</i>								
BPE w/ Beam Seach	3.05	2.50	2.18	2.04	2.95	4.28	2.45	2.50
BPE w/ SubMerge	2.92	2.48	2.15	2.10	2.88	4.00	2.34	2.39
BPE-dropout w/ Beam Seach	4.30	3.29	3.19	2.95	4.76	8.01	4.13	4.54
BPE-dropout w/ SubMerge	3.56	2.92	2.70	2.53	3.85	6.06	3.35	3.68
<i>BeamSize=3</i>								
BPE w/ Beam Seach	2.98	2.45	2.15	2.01	2.88	4.09	2.41	2.47
BPE w/ SubMerge	2.84	2.44	2.12	2.05	2.79	3.83	2.30	2.36
BPE-dropout w/ Beam Seach	4.31	3.26	3.19	2.92	4.77	7.82	4.12	4.51
BPE-dropout w/ SubMerge	3.25	2.67	2.46	2.32	4.29	5.28	3.02	3.34
<i>BeamSize=4</i>								
BPE w/ Beam Seach	2.95	2.43	2.14	2.00	2.85	4.03	2.39	2.45
BPE w/ SubMerge	2.79	2.42	2.11	1.98	2.76	3.79	2.40	2.36
BPE-dropout w/ Beam Seach	4.31	3.26	3.19	2.91	4.77	7.74	4.11	4.50
BPE-dropout w/ SubMerge	3.09	2.56	2.34	2.21	3.25	4.94	2.86	3.15
<i>BeamSize=5</i>								
BPE w/ Beam Seach	2.93	2.42	2.13	1.99	2.83	3.98	2.38	2.45
BPE w/ SubMerge	2.77	2.41	2.10	2.02	2.75	3.75	2.35	2.35
BPE-dropout w/ Beam Seach	4.26	3.24	3.18	2.91	4.80	7.68	4.10	4.49
BPE-dropout w/ SubMerge	2.97	2.52	2.26	2.15	3.48	4.68	2.75	3.03
<i>BeamSize=6</i>								
BPE w/ Beam Seach	2.92	2.41	2.13	1.99	2.82	3.97	2.38	2.44
BPE w/ SubMerge	2.76	2.40	2.10	1.97	2.75	3.75	2.35	2.35
BPE-dropout w/ Beam Seach	4.27	3.23	3.18	2.91	4.76	7.64	4.10	4.49
BPE-dropout w/ SubMerge	2.88	2.46	2.21	2.09	3.09	4.51	2.68	2.95
<i>BeamSize=8</i>								
BPE w/ Beam Seach	2.90	2.40	2.12	2.00	2.82	3.94	2.38	2.44
BPE w/ SubMerge	2.74	2.39	2.09	2.01	2.74	3.72	2.32	2.34
BPE-dropout w/ Beam Seach	4.23	3.22	3.19	2.90	4.76	7.56	4.08	4.48
BPE-dropout w/ SubMerge	2.82	2.43	2.16	2.05	2.98	4.35	2.63	2.87
<i>BeamSize=10</i>								
BPE w/ Beam Seach	2.89	2.40	2.12	2.00	2.82	3.93	2.38	2.44
BPE w/ SubMerge	2.74	2.39	2.09	1.98	2.73	3.72	2.32	2.34
BPE-dropout w/ Beam Seach	4.20	3.22	3.17	2.89	4.69	7.51	4.08	4.46
BPE-dropout w/ SubMerge	2.79	2.41	2.13	2.03	3.03	4.22	2.58	2.82

Table 8: Word perplexity of the generated hypothesis using different beam sizes.

	IWSLT'15 Vi→En	IWSLT'15 En→Vi	WMT'16 Ro→En	WMT'16 En→Ro	WMT'15 Fi→En	WMT'15 En→Fi	WMT'14 De→En	WMT'14 En→De
<i>BeamSize=1</i>								
BPE w/ Beam Seach	23.68	24.44	31.34	32.53	16.61	14.10	29.34	25.16
BPE w/ SubMerge	24.22	24.18	31.08	32.14	16.55	14.05	26.45	24.91
BPE-dropout w/ Beam Seach	29.28	28.52	34.45	34.87	18.21	16.09	29.08	24.46
BPE-dropout w/ SubMerge	29.13	28.56	34.10	34.53	18.31	16.28	-	24.31
<i>BeamSize=2</i>								
BPE w/ Beam Seach	23.98	24.92	31.82	32.79	17.06	14.83	30.05	25.63
BPE w/ SubMerge	25.40	24.80	31.45	32.62	17.10	14.80	30.19	25.59
BPE-dropout w/ Beam Seach	29.87	29.21	35.02	35.25	18.64	16.70	29.48	24.73
BPE-dropout w/ SubMerge	30.01	29.16	34.56	35.19	18.80	16.30	29.53	24.72
<i>BeamSize=3</i>								
BPE w/ Beam Seach	24.44	24.92	32.03	33.02	16.89	15.30	30.25	25.84
BPE w/ SubMerge	25.47	24.94	31.66	33.00	16.97	14.84	30.23	25.75
BPE-dropout w/ Beam Seach	29.77	29.34	35.06	35.55	18.77	16.98	29.58	24.89
BPE-dropout w/ SubMerge	29.64	29.20	34.45	35.36	18.26	16.49	29.51	24.84
<i>BeamSize=4</i>								
BPE w/ Beam Seach	24.34	25.09	32.05	32.98	17.08	15.30	30.18	25.88
BPE w/ SubMerge	25.63	24.86	31.70	32.85	16.94	15.06	30.04	25.71
BPE-dropout w/ Beam Seach	29.65	29.40	34.96	35.44	18.80	16.95	29.75	24.79
BPE-dropout w/ SubMerge	30.03	29.61	34.75	35.39	18.87	16.64	29.70	24.94
<i>BeamSize=5</i>								
BPE w/ Beam Seach	24.38	25.02	32.02	32.95	17.05	15.26	30.13	25.80
BPE w/ SubMerge	25.79	24.93	31.75	33.00	17.07	14.89	30.08	25.67
BPE-dropout w/ Beam Seach	29.36	29.44	34.99	35.55	18.97	17.14	29.69	24.82
BPE-dropout w/ SubMerge	29.50	28.67	34.48	35.43	18.37	16.74	29.61	24.87
<i>BeamSize=6</i>								
BPE w/ Beam Seach	24.45	25.07	31.99	32.88	17.02	15.28	30.11	25.78
BPE w/ SubMerge	25.58	24.86	31.62	32.85	17.04	14.77	30.07	25.67
BPE-dropout w/ Beam Seach	29.35	29.53	34.96	35.47	19.05	17.26	29.61	24.81
BPE-dropout w/ SubMerge	29.58	29.18	34.46	35.32	14.41	16.99	29.50	24.91
<i>BeamSize=8</i>								
BPE w/ Beam Seach	24.57	24.81	31.86	32.37	17.14	15.23	30.06	25.68
BPE w/ SubMerge	25.86	24.83	31.72	32.81	16.78	14.89	29.93	25.63
BPE-dropout w/ Beam Seach	29.33	29.52	34.99	35.58	18.95	17.21	29.54	24.73
BPE-dropout w/ SubMerge	29.73	29.52	34.76	35.27	18.88	17.06	29.46	24.87
<i>BeamSize=10</i>								
BPE w/ Beam Seach	24.77	24.93	31.78	32.05	17.09	15.24	30.09	25.75
BPE w/ SubMerge	25.93	24.66	31.63	32.85	16.69	14.97	29.90	28.16
BPE-dropout w/ Beam Seach	29.32	29.41	34.97	35.59	19.08	17.32	29.58	24.80
BPE-dropout w/ SubMerge	29.39	28.92	34.45	35.39	18.88	16.92	29.28	24.68

Table 9: BLEU scores on test sets using different beam sizes.

FAME-MT Dataset: Formality Awareness Made Easy for Machine Translation Purposes

Dawid Wiśniewski^{1,2}, Zofia Rostek¹, Artur Nowakowski^{1,3}

¹ Lanigo, Poznań, Poland

² Faculty of Computing and Telecommunications, Poznań University of Technology, Poland

³ Faculty of Mathematics and Computer Science, Adam Mickiewicz University, Poznań, Poland
name.surname@lanigo.com

Abstract

People use language for various purposes. Apart from sharing information, individuals may use it to express emotions or to show respect for another person. In this paper, we focus on the formality level of machine-generated translations and present **FAME-MT** – a dataset consisting of 11.2 million translations between 15 European source languages and 8 European target languages classified to formal and informal classes according to target sentence formality. This dataset can be used to fine-tune machine translation models to ensure a given formality level for each European target language considered. We describe the dataset creation procedure, the analysis of the dataset’s quality showing that **FAME-MT** is a reliable source of language register information, and we present a publicly available proof-of-concept machine translation model that uses the dataset to steer the formality level of the translation. Currently, it is the largest dataset of formality annotations, with examples expressed in 112 European language pairs. The dataset is published online¹.

1 Introduction

Motivation Different situations require using different language depending on whether it is a formal meeting or a casual talk with friends. Fre-

quently, when speaking to an older or distinguished person, e.g., an owner of a company or a university professor, we use appropriate language forms to show respect. However, machine translation models may struggle with choosing an appropriate language form due to lack of context – often, we only have one sentence to translate, there are cultural differences between source and target language speakers, and models may be trained on parallel corpora, which do not focus on formality aspect of language enough.

For this reason, it is important to find ways of enforcing machine translation models to use a required formality level in translated sequences. While there are methods that can be used to achieve this goal, the existing datasets are scarce, focusing either on formality classification for a few selected languages (e.g., German or English) or providing formality-annotated examples of translations between pairs of languages, which are limited in size and include a narrow selection of target languages with English as the source language. What we found missing is a large-scale dataset providing formality-annotated examples of translations for a much broader set of languages, enabling easy fine-tuning of pre-trained machine translation models not only for directions including English but also other European languages.

Contribution In this paper, we present **FAME-MT**, the biggest to date dataset of translation pairs annotated with formality level. This dataset introduces 100,000 annotated translation examples for each of the 112 European language pairs considered. With 8 target languages (English, German, French, Italian, Dutch, Polish, Portuguese, and Spanish) classified according to formality and 15 European languages considered as the source lan-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹<https://github.com/lanigo-public/fame-mt/>

guages, the dataset introduces an 18-fold increase of the European language-pair coverage over the most diverse dataset **CoCoA-MT** (Nadejde et al., 2022) and over 100-fold increase of the dataset size considering the biggest formality-annotated datasets available (Rao and Tetreault, 2018).

Providing translation examples for a wide selection of source and target languages, **FAME-MT** allows for simple fine-tuning of machine translation models for the most popular European language pairs.

In this paper, we discuss the dataset creation process and analyze the data quality using various metrics to prove the usefulness of **FAME-MT**. We show how proof-of-concept models for formality-aware machine translation can be trained using Marian (Junczys-Dowmunt et al., 2018), provide examples of outputs of those models, and publish models online along with the dataset.

Research questions We formulated the following research questions to be answered in this paper:

- **RQ1:** Is it possible to create a good quality large-scale dataset for formality-aware machine translation automatically based on available resources?
- **RQ2:** Is 100,000 translation examples for a given language pair enough to fine-tune a pre-trained machine translation model to become formality-aware?
- **RQ3:** Considering translation pairs coming from **MTData** (Gowda et al., 2021), are formal sentences always translated into formal ones, and informal sentences translated into informal ones?

2 Related work

The idea of identifying formality level in texts is a widely analyzed area of (socio)linguistics (Biber and Conrad, 2019). In linguistics, there are 5 main language registers defined that can be used in particular situations. These are: frozen, formal, consultative, casual (informal), and intimate registers. Out of them, the formal and informal ones are most often analyzed in the context of machine translation.

Existing formality datasets There are several datasets proposed that help incorporate formality

awareness in various NLP tasks. The biggest one, Grammarly’s Yahoo Answers Formality Corpus (or **GYAFC** corpus for short (Rao and Tetreault, 2018)), introduces 110,000 pairs of formal and informal sentences in English. The dataset is based on Yahoo Answers L6 corpus² and is proposed to be used for style transfer. The authors proposed an LSTM-based baseline showing that style transfer can be achieved using the machine translation approach and propose metrics for automatic evaluation of style-transfer models (e.g., fluency, meaning preservation).

Another dataset, consisting of 6,574 English sentences annotated with formality level, is proposed by Pavlick and Tetreault (2016). This work is aimed at the analysis of the formality level in English, considering humans’ perceptions of formality in four different genres.

There are also datasets that do not focus on English, for example, a dataset consisting of 3,000 German sentences annotated by human experts on a continuous scale using comparative judgments (Eder et al., 2023). Each annotator was presented with several examples, and their goal was to rank the sentences according to formality level. This dataset covers a set of 12 diverse data sources (e.g., Twitter, Reddit, Wikipedia, or Springer Open Science articles).

Apart from the datasets that focus on a single language, there are several datasets with multilingual data. One of them is **XFORMAL** (Briakou et al., 2021). Similarly to **GYAFC**, the goal of **XFORMAL** is to provide a benchmark for style transfer by introducing pairs of informal sentences and their formal counterparts. The dataset provides examples in three languages: French, Italian, and Brazilian-Portuguese, and introduces 1,000 human-annotated examples for each of these languages. However, access to both **GYAFC** and **XFORMAL** is limited as it requires access approval as described on the website³.

A publicly available dataset called **CoCoA-MT** (Nadejde et al., 2022) provides a set of triples consisting of English sentences translated into formal and informal versions for each of the supported languages, namely, German, Spanish, Italian, French, Dutch, Portuguese, Japanese, and Hindi. The dataset provides around 1,000 examples divided into train and test sets for each tar-

²<https://webscope.sandbox.yahoo.com/catalog.php?datatype=1>

³<https://github.com/Elbria/xformal-FoST/>

get language. What distinguishes **CoCoA-MT** is the presence of phrase annotations that tell which phrases make a given example formal or informal. Even though **CoCoA-MT** provides translations between language pairs, they are always translations from English to one of the supported languages.

In general, although there are datasets focusing on formality, the large ones (**GYAFC**) cover only English or German, while those including other languages are relatively small (**XFORMAL**) and provide translations only from English to a given language (**CoCoA-MT**). No large-scale multi-language dataset exists up to date.

Formality control methods The tasks of style transfer and machine translation can be controlled using various methods. Here, we discuss the most relevant ones in the context of formality control.

The most straightforward method of controlling formality injects a special token in a source sequence, which tells us what level of formality should be achieved on the target side. This method was utilized in various approaches, e.g.: controlling honorifics in English to Japanese translation (Feely et al., 2019) using one of `{informal, polite, formal}` tags injected to the input of the Transformer model, controlling formality of French to English translations using one of predefined formality levels (Niu et al., 2017) `{low, neutral, high}`, controlling formality level by injecting a special token representing one of `{informal, formal}` classes to control output formality presented as part of **CoCoA-MT** evaluation (Nadejde et al., 2022), controlling politeness of the output text using predefined Latin tokens: `{vos, tu}` (Sennrich et al., 2016), or controlling formality by attaching a special token to both input and output sequences for better control of the output (Niu and Carpuat, 2020).

Alternatively, one can add a special embedding vector for each token (Schioppa et al., 2021) to represent the desired formality level, craft prompts (Garcia and Firat, 2022) for multilingual T5 (Xue et al., 2021) models, or use Bayesian factorization for constrained output generation (Yang and Klein, 2021).

3 Dataset

The **FAME-MT** dataset creation process was divided into three steps.

Step 1: Input data selection The aim of the **FAME-MT** project was to provide translations between pairs of languages that are annotated with formality information. For this purpose, the first step was to identify parallel corpora for languages of interest. We considered eight target languages, i.e., languages into which we want to translate input statements. These are: *English (EN)*, *German (DE)*, *French (FR)*, *Italian (IT)*, *Dutch (NL)*, *Polish (PL)*, *Portuguese (PT)*, and *Spanish (ES)*. For each target language, our goal was to support a wide selection of source languages, i.e., those from which we translate into the target language.

We selected 15 source languages: *Czech (CS)*, *Danish (DA)*, *German (DE)*, *English (EN)*, *Spanish (ES)*, *French (FR)*, *Italian (IT)*, *Dutch (NL)*, *Norwegian* – including Norwegian Bokmål (*NO + NB*), *Polish (PL)*, *Portuguese (PT)*, *Russian (RU)*, *Slovak (SK)*, *Swedish (SV)*, and *Ukrainian (UK)*. This selection of source and target languages resulted in 112 potential language pairs, for each of which we needed a parallel corpus of translations between the given language pair. To collect such a dataset, we used **MTData** tool that provides an access to machine translation dataset collections, e.g., OPUS (Tiedemann, 2016), gives an access to popular datasets such as Europarl (Koehn, 2005), or Paracrawl (Bañón et al., 2020), and covers every language pair considered. For each language pair, we acquired a corpus of translation examples from **MTData** and applied the following post-processing to increase the quality of the collected data: we rejected documents with more than 15% of characters being digits, having less than 5 characters in any sentence, having any token longer than 28 characters, having any sentence longer than 500 characters. We also rejected those with the number of tokens in any sentence higher than 100. Moreover, we applied FastText’s LID-201 model (Burchell et al., 2023) to verify whether the source and target sentences are indeed expressed in the correct language (the expected language should have probability score of at least 10%).

For language pairs where English is present as a source or target language, we used Bicleaner AI package, which estimates the likelihood of a pair of sentences being mutual translations (Zaragoza-Bernabeu et al., 2022). We set the bicleaner-score threshold to 50%, rejecting everything below this score. Bicleaner cannot be applied to other languages as for now because it does not provide

open-source models for such language pairs.

Step 2: Formality classification The next step involved extracting two subsets for each language pair: those with formal and informal translations, respectively. For this reason, for each target language, we searched for formality classifiers or golden standard annotations that can be used to train a classifier. As discussed in Section 6, we only need to classify the target language according to its formality level. There were three scenarios to address:

- Formality classifier available: For English, there is a publicly available classifier available online ⁴, the quality of which is proven by an accompanying research paper (Babakov et al., 2023). As the classifier was pre-trained on the biggest dataset available (**GYAFC**) and produces the probability of a given example to be formal, we selected this classifier to process those pairs of sentences where the target language is English.
- No formality classifier, but golden standard dataset available: For German, French, Italian, Spanish, Portuguese, and Dutch, no formality classifier accompanied by a research paper could be found. For this reason, we decided to train a classifier using **CoCoA-MT** dataset, as it provides pairs of formal and informal forms of sentences for each of the aforementioned target languages.
- No formality classifier nor dataset available: For Polish, we did not even have a golden standard dataset to train a classifier on. For that language, we created a hold-out subset of sentences in Polish downloaded using **MT-Data**, and then the set was annotated by a group of six native speakers. This way, we collected examples of formal and informal sentences, which were of roughly the same size as golden standard examples provided in **CoCoA-MT** and were split into train and test sets to mimic the structure of **CoCoA-MT**.

We decided to include the Polish dataset in examples collected from **CoCoA-MT** and tried to fine-tune a single multilingual classifier to capture inter-lingual relations.

⁴<https://huggingface.co/s-nlp/roberta-base-formality-ranker>

However, early experiments with models described later in this Section showed that fine-tuning a pre-trained model using **CoCoA-MT** combined with Polish leads to extreme probabilities assigned to most sentences collected from **MTData**. As the dataset contains only formal and informal examples, classifiers learn to treat every sentence as either formal or informal, while in many cases (especially sentences that are not related to an interlocutor) they are neither formal nor informal.

An analysis of annotations from **CoCoA-MT** showed that formality and informality are frequently expressed using appropriate personal pronouns. For this reason, we decided to generate an additional neutral dataset for each language generated as a random hold-out sample of target sentences from **MTData** that does not contain any phrase marked in **CoCoA-MT** as either formal or informal. For Polish, we continued the annotation task, asking natives to identify neutral examples in the hold-out dataset.

We fine-tuned several language models verifying how big a neutral sample size should be to maximize the scores. At each verification step, we changed only the size of the neutral set in the train set, while preserving a constant set of 600 neutral examples in the test set. We evaluated accuracy score for mDeberta-v3-base (He et al., 2021), XLM-RoBERTa-uncased (Conneau et al., 2019), and BERT-base-multilingual-uncased (Devlin et al., 2019). The results, summarized in Table 1 show that the best average accuracy score is achieved using mDeberta-v3-base with a neutral training sample size = 500. This comes in line with the trainset size of original **CoCoA-MT**, as it uses 400 formal and 400 informal examples for each language pair, so including a neutral sample of this size results in an approximately balanced training dataset.

Diving deeper into per-language scores, as presented in Table 2, we observe that all **CoCoA-MT** languages obtain very high accuracy scores. This may be due to the fact that frequently examples in **CoCoA-MT** distinguish formality based on personal pronouns. For Polish, the scores are lower as personal pronouns are often dropped and the form of a verb may be used to express a linguistic person implicitly. However, even for Polish, the scores are much higher than random guesses. For English, we report the scores for **GYAFC** obtained from a research paper that a given model accom-

NSS	mDeberta	XLM-RoBERTa	mBERT
100	0.9304	0.9301	0.9297
200	0.9475	0.9437	0.9395
300	0.9519	0.9465	0.9454
400	0.9531	0.9511	0.9455
500	0.9552	0.9523	0.9499
600	0.9551	0.9534	0.9492

Table 1: Neutral sample size (NSS) for training vs. average accuracy score (over all languages) of a given model. mDeberta stands for mDeberta-v3-base, XLM-RoBERTa stands for XLM-RoBERTa-uncased, and mBERT stands for BERT base multilingual uncased.

panies (Babakov et al., 2023).

Step 3: FAME-MT compilation The dataset compilation process was performed as follows: For each language *src* from source languages set and each target language *tgt*, we classified all targets among translations between *src* and *tgt*. When the target language was English, we used a model trained on **GYAFC** (Babakov et al., 2023). For other languages, we used mDeberta-v3-base, which we fine-tuned in the previous step. Since mDeberta-v3-base returns probability distribution over three possible classes: formal, informal, and neutral, we chose the class with the highest probability as the final model decision. However, since the model trained on **GYAFC** returns only the probabilities of formal and informal classes (that always sum up to one), we split the probability range into three equal parts, treating examples assigned with formal class probability in ranges: $< 0, \frac{1}{3} >$, $(\frac{1}{3}, \frac{2}{3} >$, $(\frac{2}{3}, 1 >$ as informal, neutral, and formal, respectively.

The classification process continued until we reached 50,000 informal and 50,000 formal examples for each language pair. Finally, for each language pair, we stored translations where the target sentence was considered formal or informal into separate files.

The dataset is published online⁵ along with the scripts used for analysis, formality classifiers, and MT models. It has the following structure: Each target language has its own folder assigned, and inside these folders, there is a separate folder for each source and target language combination. In those folders, there is a pair of files: *informal.tsv* and *formal.tsv* each

⁵<https://github.com/lanigo-public/fame-mt/>

providing 50,000 translations between source and target languages, where the target was considered informal and formal, respectively.

4 Explorative analysis

To verify the quality of the dataset, we selected several approaches that would show key characteristics of the dataset.

Average length of sentences One of the basic metrics that can describe the difference between formal and informal language is the average sentence length in each of the categories. We may expect that formal sentences may be longer than informal, which are frequently used to share knowledge quickly. Figure 1 proves that for each target language, the formal translations are much longer than informal ones in the collected dataset.

Per-category key tokens To check which tokens characterize a given class, we aimed at understanding classifier decisions and feature importance. As the classifiers we used in the dataset creation process were too big to be handled by the popular explainable AI framework LIME (Ribeiro et al., 2016), we decided to fit a lightweight linear model to **FAME-MT** and interpret its features to understand the difference between classes.

Having a set of sentences classified as formal and informal for each target language, we applied the TF-IDF vectorizer to a class-balanced sample of target sentences to identify the 200,000 most important tokens among target texts. Then, we used those tokens as features for a logistic regression model with the intercept value set to 0. This way, we ensure that the coefficients assigned by the logistic regression model correlate with a given class. As we modeled formal sentences as a positive class and informal ones as negative, tokens that are assigned with high positive coefficient values are strongly related to the formal class, and those associated with the highly negative values are strongly related to the informal class. The coefficients selected during training were then sorted for each language, and the top 100 tokens and bottom 100 were selected as those describing formal and informal categories, respectively.

Table 3 presents the most important features for each category and each language considered. As can be seen, for languages with formal *you*, the most important words are personal pronouns (e.g., formal *Sie* vs. informal *du* in German or formal

Language	German	French	Italian	Dutch	Polish	Portuguese	English
Accuracy	0.9928	0.9926	0.9772	0.9962	0.7861	0.9789	0.9 (Babakov et al., 2023)

Table 2: Classification scores for each language using mDeberta-v3-base and neutral sample size = 500 in comparison to English classifier trained on **GYAFC**. The average score of mDeberta-v3-base is equal to 0.9552.

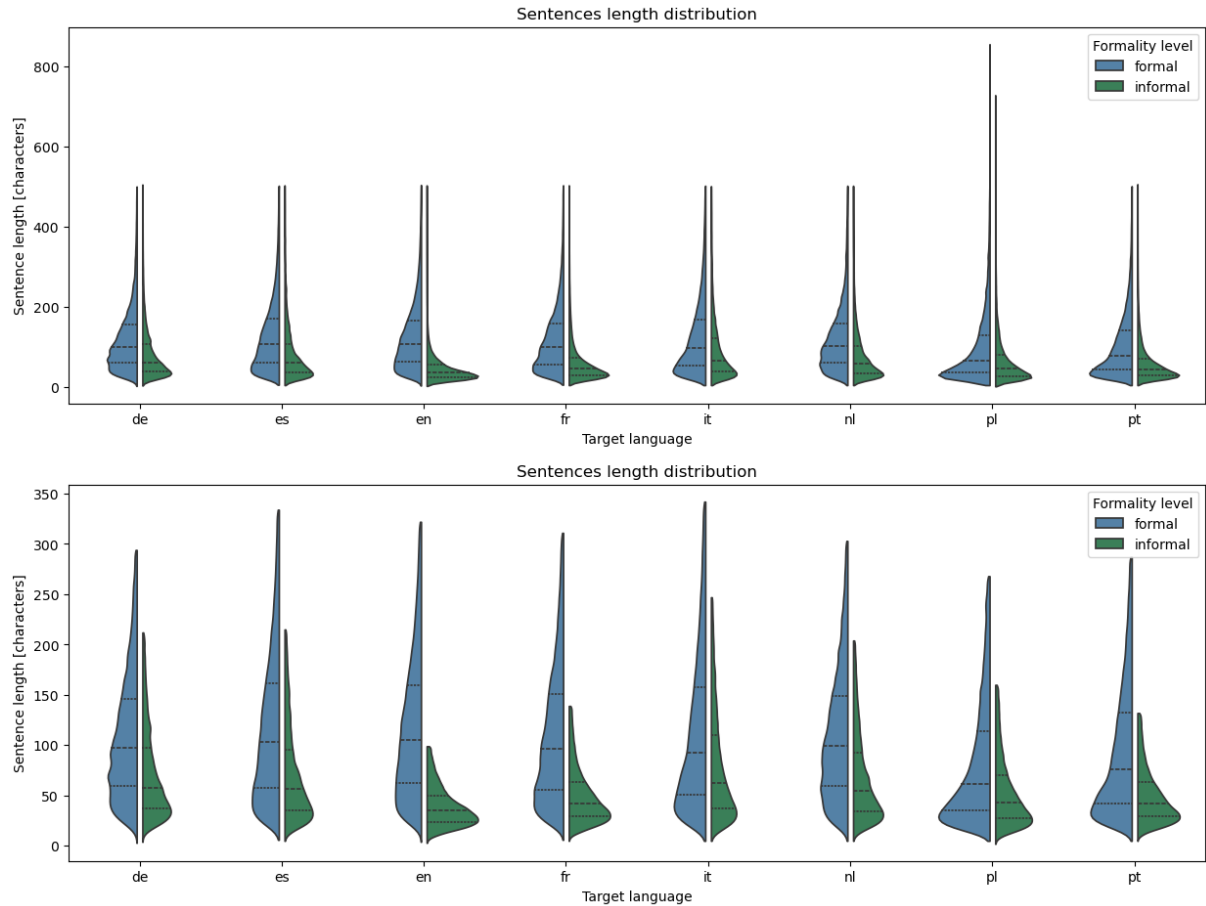


Figure 1: Violin plots representing the distributions of sentence lengths interpreted as the number of characters. The upper figure represents the distributions calculated over the original dataset. As it shows that there are some outliers with big values, we provide the lower figure generated over a subset of texts whose lengths are between $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$ ($Q1$ =first quartile, $Q3$ =third quartile, IQR =inter-quartile range) to focus more on the most common scenarios.

usted and informal *tu* in Spanish). For languages, where personal pronouns may be omitted in informal scenarios, the grammar form of verbs may indicate informality (e.g., *jesteś* (informal *you are*) in Polish). For English, where there is no formal *you*, more sophisticated words are observed as the most important formal tokens, and contractions or swear words are found among the most characteristic informal tokens.

Textual complexity We used textacy⁶ to measure the complexity of target sentences classified as formal and informal. The motivation for this step was the intuition that informal language may be easier to comprehend while formal texts should

be harder. For English, we used the automated readability index function, which calculates the relation between the number of characters in a given sequence related to the number of words and sentences and is calculated using the following formula: $4.71 \cdot \text{count}_{chars} / \text{count}_{words} + 0.5 \cdot \text{count}_{words} / \text{count}_{sents} - 21.43$. We calculated the score for all language pairs with English as a target and calculated the mean value. We obtained 11.54 for formal sentences and 4.28 for informal ones, which means that formal sentences are more complex than informal ones. Considering each language pair with the target set to English, formal translations are more complex every time. The biggest difference between formal and infor-

⁶<https://pypi.org/project/textacy/>

Language	formal words	informal words
German (DE)	sie, ihre, ihr, ihnen, ihres, ...	du, deine, dir, dich, dein, ...
English (EN)	distinctive, relations, moreover, obtain, refers, ...	gotta, f****g, gonna, ain, wanna, ...
Spanish (ES)	su, sus, usted, le, está, ...	tu, te, tus, estás, quieres, ...
French (FR)	vous, votre, vos, pouvez, avez, ...	tu, toi, te, ton, ta, ...
Italian (IT)	sua, suo, suoi, lei, le, voi, ...	ti, hai, tuo, tua, tuoi, sai, tu, ...
Dutch (NL)	uw, kunt, wilt, bent, hebt ...	je, jij, jullie, jouw, jou, ...
Polish (PL)	pan, pani, wam, państwa, państwo, ...	jesteś, ci, możesz, myślisz, musisz, ...
Portuguese (PT)	ocê, sua, seu, suas, seus, ...	te, teu, tua, tu, tens, ...

Table 3: Words with the strongest relation to classes selected by analyzing feature coefficients in logistic regression model fitted on FAME-MT.

mal sentence scores is seen in the case of French-English (formal: 12.856, informal: 2.788) and the smallest in the case of Czech and English (formal: 11.175, informal: 7.005). We analyzed the Flesch score, which indicates the reading ease for sentences collected. We sampled 5,000 sentences for each source - target language pair, where target sentence was among languages supported by textacy (German, English, French, Italian, Spanish, Portuguese, Dutch). We present the full analysis in Table 10 provided in Appendix A. For each language pair, the informal texts are scored higher than formal ones, which indicates that informal texts are easier to read. This follows the general intuition. Averaging over all source languages, the average Flesch scores for formal and informal texts are (F and I represent formal and informal, respectively): German: 51.99 F and 72.57 I, English: 63.19 F & 90.89 I, Spanish: 68.93 F & 80.45 I, French: 78.33 F & 99.69 I, Italian: 77.65 F & 81.4 I, Dutch: 57.61 F & 76.32 I, Portuguese: 66.9 F & 76.32 I.

Source vs. target sentence formality To explore the relation between the formality level of source sentence and its translation in target language, we selected a subset of FAME-MT for which we had classifiers for both source and target languages. As these are eight languages, from all 112 language pairs, we selected 56 of them (8 times 8 possible language pairs minus 8 pairs where both source and target are the same). Then, for each target language considered, we analyzed formal translations and informal translations separately: for each formal translation, we used appropriate classifier to verify whether a given source sentence is also formal and for each informal translation, we used appropriate classifier to verify whether a given source sentence is also informal. Then, we calculated the percentage of cases where the agreement was observed. As presented

in Table 4, on average, 38.81% of examples for which target sentence was considered formal had the source sentence classified as formal. Analogously, 40.28% of examples classified as informal had the source sentence classified as informal. These results show that the formality level of the source does not determine the formality level of the target text. However, there is some correlation as the scores are higher than a random choice, which would achieve 33% (assigning a random class from: formal, informal, neutral). Diving into language pairs as presented in Table 11 in Appendix A and considering translations that are formal in target language, the strongest relation is between French sources and Dutch targets (82.8% of formal cases in Dutch are also formal in French), and the weakest between Polish sources and English targets (5.6% of formal cases in English are also formal in Polish). Considering translations that are informal in target language, the strongest relation is between Polish sources and French targets (73.4% of informal cases in French are also informal in Polish), and the weakest between German sources and English targets (5.6% of informal cases in English are also informal in German).

5 Application to Machine Translation

This section explores the application of the FAME-MT dataset to enable control over the formality tone of translated text. We conducted experiments on two translation directions: English-German and English-Polish, leveraging pre-trained models from the OPUS collection (Tiedemann, 2020) for fine-tuning.

Fine-tuning pre-trained MT models To facilitate the fine-tuning process, we augmented the existing model vocabulary with two new tokens: <FORMAL> and <INFORMAL>. Additionally, we expanded the existing model embeddings to accommodate these new vocabulary items by initial-

Language	Formal S+T	Informal S+T
German (DE)	42.82%	48.41%
English (EN)	17.71%	14.95%
Spanish (ES)	41.18%	40.82%
French (FR)	45.51%	54.4%
Italian (IT)	37.91%	37.36%
Dutch (NL)	55.56%	39.38%
Polish (PL)	37.16%	30.87%
Portuguese (PT)	35.67%	56.04%
Average	38.81%	40.28%

Table 4: Percentage of examples where formal source co-occurs with formal target (Formal S+T), and percentage of examples, where informal source co-occurs with informal target (Informal S+T). Formality of source and targets was determined using classifiers. The row *Average* represents the micro-average over all languages.

izing them with zeros.

To maintain consistency with the baseline models, we reused their Transformer-based (Vaswani et al., 2017) architecture and training hyperparameters, except for batch size and learning rate. The English-German direction utilizes the Transformer-big architecture, while the English-Polish direction uses the Transformer-base architecture due to the lack of a larger model in the OPUS collection for this translation direction. Both fine-tuned models employed a small batch size of 5,000 target tokens and a learning rate of $1e-4$. As the original models were trained using the Marian (Junczys-Dowmunt et al., 2018) framework, we utilized it for fine-tuning as well.

We pre-processed the **FAME-MT** dataset with pre-trained SentencePiece (Kudo and Richardson, 2018) tokenizers included in the model package, and then prepended `<FORMAL>` and `<INFORMAL>` tags to the formal and informal parts of the source input, respectively. To guarantee a diverse validation set during training, we randomly sampled 500 sentences from each formality category and added 500 random neutral samples not using personal pronouns, resulting in a total of 1500 validation samples per translation direction. The remaining data was used for training.

Automatic evaluation and contrastive examples

Three decoding modes were evaluated for the fine-tuned models: standard, formal tone, and informal tone. Their performance was compared to the baseline OPUS model on the Flores (devtest) (Costa-jussà et al., 2022) and NTREX (Ferdemann et al., 2022) datasets, using BLEU (Pa-

pineni et al., 2002), chrF (Popović, 2015), and COMET⁷ (Rei et al., 2020) metrics for evaluation. The chrF and BLEU results were computed with sacreBLEU⁸⁹ (Post, 2018). The evaluation results are presented in Tables 5 and 6 for Flores and NTREX datasets, respectively.

The results for English-German translation were encouraging. Fine-tuning on **FAME-MT** did not negatively impact overall translation quality, and automatic metrics even suggested a slight improvement. However, the English-Polish model exhibited a slight decrease in quality, potentially due to the smaller model size being more susceptible to the impact of embedding extension.

Tables 7 and 8 showcase various translations of the same sentence with different formality levels for each fine-tuned model. Interestingly, in languages like Polish, where gender neutrality is not present in singular pronouns, formal translations can introduce gender bias. For example, English *you* can be translated as either *Pan* (masculine) or *Pani* (feminine). In such cases, sentence-level machine translation might struggle to choose the correct form without additional context. As expected, formal translations generally tend to be more literal compared to their informal counterparts.

Released models To promote further research and facilitate accessibility, we have made the fine-tuned models publicly available as open-source resources¹⁰.

This work demonstrates the potential of incorporating formality control datasets into machine translation pipelines. While further investigation is needed to refine the approach for different model sizes and language pairs, the results show promise for producing translations that accurately reflect the desired level of formality. The open-sourced models will enable researchers to build upon this work and explore applications in diverse scenarios.

6 Discussion

As we have shown in the previous section, the **FAME-MT** dataset can make pre-trained models formality-aware. Also, metrics used to analyze the dataset show that the characteristics of the dataset

⁷[wmt22-comet-da](https://github.com/chrF/comet-da) COMET model was used

⁸BLEU signature: nrefs:1lcase:mixedleff:noltok:13a lsmooth:explversion:2.3.1

⁹chrF signature: nrefs:1lcase:mixedleff:yeslnc:6lnw:0 lspace:nolversion:2.3.1

¹⁰<https://github.com/lanigo-public/fame-mt/>

Model		Flores					
		English → German			English → Polish		
		COMET	chrF	BLEU	COMET	chrF	BLEU
OPUS-Finetuned-50k	Standard	0.8687	66.14	39.98	0.8566	50.54	19.87
	Formal	0.8656	66.18	40.02	0.8523	50.47	19.70
	Informal	0.8687	65.85	39.60	0.8539	50.48	19.84
OPUS (baseline)	Standard	0.8675	66.20	39.90	0.8624	50.85	20.15

Table 5: Automatic evaluation results on the Flores dataset before and after fine-tuning the model on the FAME-MT dataset.

Model		NTREX					
		English → German			English → Polish		
		COMET	chrF	BLEU	COMET	chrF	BLEU
OPUS-Finetuned-50k	Standard	0.8295	60.25	32.23	0.8097	51.45	23.03
	Formal	0.8256	60.21	31.98	0.8049	51.32	22.83
	Informal	0.8328	60.15	31.95	0.8047	51.22	22.77
OPUS (baseline)	Standard	0.8251	59.91	31.64	0.8125	51.89	23.41

Table 6: Automatic evaluation results on the NTREX dataset before and after fine-tuning the model on the FAME-MT dataset.

English	Polish (Formal)	Polish (Informal)
You have to tip your cap.	Musi pan przechylić czapkę.	Musisz przechylić czapkę.
They think you're sad and will be pleased because they got to you.	Sądzą, że jest pan smutny i będą zadowoleni, bo do pana doszli.	Sądzą, że jesteś smutny i będą zadowoleni, bo cię dorwali.
They don't know you're furious.	Nie wiedzą, że jest pani wściekła.	Nie wiedzą, że jesteś wściekła.
But on Saturday, North Korean Foreign Minister Ri Yong-ho blamed US sanctions for the lack of progress since then.	Ale w sobotę, północnokoreański minister spraw zagranicznych Ri Yong-ho obarczył sankcje USA za brak postępów od tego czasu.	Ale w sobotę, północnokoreański minister spraw zagranicznych Ri Yong-ho obwiniał USA o sankcje za brak postępów od tego czasu.
Ring also settled a lawsuit with competing security company, the ADT Corporation.	Ring rozstrzygnął również sprawę z konkurencyjnym koncernem ochroniarskim, korporacją ADT.	Ring rozstrzygnął również sprawę z konkurencyjną firmą ochroniarską, firmą ADT Corporation.

Table 7: Examples of different machine translation results with formal/informal tone in English → Polish translation.

English	German (Formal)	German (Informal)
You just have to pay them the right amount of respect, he said..	Du musst ihnen nur den richtigen Respekt erweisen, sagte er.	Sie müssen ihnen nur die richtige Menge an Respekt zahlen, sagte er.
Sit down.	Setzen Sie sich.	Setz dich.
Would you like to go with me to the cinema next week?	Möchtest du nächste Woche mit mir ins Kino gehen?	Möchten Sie nächste Woche mit mir ins Kino gehen?
I'll be with you soon.	Ich werde bald bei dir sein.	Ich bin bald bei Ihnen.
I would like a word with your boss.	Ich hätte gerne ein Wort mit deinem Chef.	Ich möchte ein Wort mit Ihrem Chef.

Table 8: Examples of different machine translation results with formal/informal tone in English → German translation.

are consistent with intuition: formal documents are harder to read, they tend to be longer or use more formal words.

An interesting observation is that even though the **CoCoA-MT** dataset introduces a relatively small dataset of formal and informal examples per each language, the accuracy of the pre-trained mDeberta-v3-base model, which is fine-tuned on **CoCoA-MT** is very high when measured on **CoCoA-MT**'s test set with neutral samples added (0.9552). An analysis of the annotations provided in **CoCoA-MT** reveals that it focuses on expressing formality using personal pronouns. However, formality level, in general, can be expressed using various language constructs, e.g., contractions, formal greetings, slang words, and appropriate personal pronouns. On the other hand, e.g., **GYAFC** focuses more on constructs other than personal pronouns due to the lack of formal *you* in English. We think that this is the reason for the high scores observed for a small dataset (**CoCoA-MT**) and lower scores for a much bigger dataset (**GYAFC**). The focus of **CoCoA-MT** on pronouns leads to too optimistic quality estimates as compared to the evaluation provided for the English classifier.

However, since formal constructs co-occur with each other, classifiers trained using **CoCoA-MT** work well in practice. As we can see in Tables 7 and 8 – models trained with **FAME-MT** introduce subtle differences depending on the expected formality level (e.g., *hätte gerne* vs. *möchte* or *koncernem ochroniarskim* vs. *firma ochroniarska*).

The motivation for using classifiers to classify only sentences in target languages is supported by the analysis of source vs. target sentence formality. As we have shown that the formality level of the source language does not determine the formality of the target language, in **FAME-MT**, we collect sentences of various levels of formality mapped to a given target sentence formality level. This way, we can inject a token representing the desired formality level (be it *formal* or *informal*) to steer the expected formality level. Thus, having only target sentences classified, we can fine-tune machine translation models to become formality-aware, regardless of the formality level of the source sequence.

7 Addressing research questions

In this paper, we stated three research questions that can be answered now. Regarding the first

one, we show that using a set of classifiers for the parallel corporas' target languages, we created a large dataset, which is useful for training formality-aware MT models. We proved it by fine-tuning general machine translation models and utilizing metrics confirming intuitions, thus, the answer to **RQ1** is positive. Also, the experiments show that fine-tuning pre-trained OPUS models with 50,000 formal and 50,000 informal examples is enough for fine-tuning English-German and English-Polish pairs of languages. The quality of the English-German model is higher but the quality of English-Polish is also satisfactory. For this reason, we can give the positive answer to **RQ2**. Finally, the classifiers used reveal that the formality level of the source language does not determine the formality level of the target language in the datasets collected using **MTData**. This observation justifies the need to inject a special formality token to ensure a given formality level of the translation. Thus, the answer to **RQ3** is negative.

8 Conclusions

In this paper, we introduce **FAME-MT** - a dataset consisting of 11.2 million translations between 112 European language pairs, where sentences in target languages are classified as formal or informal. As the dataset is a computer-generated silver standard, we used a set of metrics to prove the good quality of the data. Moreover, proof-of-concept models fine-tuned using formality data show that the dataset can be successfully utilized in problems requiring enforcing a given formality level of the system's output. Due to its size and large number of language pairs selected, **FAME-MT** is the largest and most diverse dataset available, introducing 18 times more European language pairs than biggest existing multilingual datasets (112 language pairs vs. 6 in **CoCoA-MT**) and providing over 100 times more sentences annotated with formality level than the biggest datasets (11.2 million vs. 110,000 in **GYAFC**). We made the dataset publicly accessible, and provided all the source codes for rerunning the analysis¹¹. We hope that this dataset may help to produce better formality-aware machine translation models especially for pairs of languages that were not yet covered or underrepresented in existing datasets (e.g., Czech → French, Polish → German, or Danish → Dutch).

¹¹<https://github.com/lanigo-public/fame-mt/>

References

- Babakov, Nikolay, David Dale, Ilya Gusev, Irina Krotova, and Alexander Panchenko. 2023. Don't lose the message while paraphrasing: A study on content preserving style transfer. In Métais, Elisabeth, Farid Meziane, Vijayan Sugumaran, Warren Manning, and Stephan Reiff-Marganiec, editors, *Natural Language Processing and Information Systems*, pages 47–61, Cham. Springer Nature Switzerland.
- Bañón, Marta, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarriás, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In Jurafsky, Dan, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online, July. Association for Computational Linguistics.
- Biber, Douglas and Susan Conrad. 2019. *Register, genre, and style*. Cambridge University Press.
- Briakou, Eleftheria, Di Lu, Ke Zhang, and Joel R. Tetreault. 2021. Olá, bonjour, salve! XFORMAL: A benchmark for multilingual formality style transfer. In Toutanova, Kristina, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 3199–3216. Association for Computational Linguistics.
- Burchell, Laurie, Alexandra Birch, Nikolay Bogoychev, and Kenneth Heafield. 2023. An open dataset and model for language identification. In Rogers, Anna, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 865–879, Toronto, Canada, July. Association for Computational Linguistics.
- Conneau, Alexis, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116.
- Costa-jussà, Marta R, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*.
- Eder, Elisabeth, Ulrike Krieg-Holz, and Michael Wiegand. 2023. A question of style: A dataset for analyzing formality on different levels. In Vlachos, Andreas and Isabelle Augenstein, editors, *Findings of the Association for Computational Linguistics: EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, pages 568–581. Association for Computational Linguistics.
- Federmann, Christian, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – news test references for MT evaluation of 128 languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24, Online, nov. Association for Computational Linguistics.
- Feely, Weston, Eva Hasler, and Adrià de Gispert. 2019. Controlling japanese honorifics in english-to-japanese neural machine translation. In Nakazawa, Toshiaki, Chenchen Ding, Raj Dabre, Anoop Kunchukuttan, Nobushige Doi, Yusuke Oda, Ondrej Bojar, Shantipriya Parida, Isao Goto, and Hidayat Mino, editors, *Proceedings of the 6th Workshop on Asian Translation, WAT@EMNLP-IJCNLP 2019, Hong Kong, China, November 4, 2019*, pages 45–53. Association for Computational Linguistics.
- Garcia, Xavier and Orhan Firat. 2022. Using natural language prompts for machine translation. *CoRR*, abs/2202.11822.
- Gowda, Thamme, Zhao Zhang, Chris Mattmann, and Jonathan May. 2021. Many-to-English machine translation tools, data, and pretrained models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 306–316, Online, August. Association for Computational Linguistics.
- He, Pengcheng, Jianfeng Gao, and Weizhu Chen. 2021. DeBERTaV3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing.
- Junczys-Dowmunt, Marcin, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia, July. Association for Computational Linguistics.
- Koehn, Philipp. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of Machine Translation Summit X: Papers*, pages 79–86, Phuket, Thailand, September 13-15.

- Kudo, Taku and John Richardson. 2018. Sentence-Piece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In Blanco, Eduardo and Wei Lu, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium, November. Association for Computational Linguistics.
- Nadejde, Maria, Anna Currey, Benjamin Hsu, Xing Niu, Marcello Federico, and Georgiana Dinu. 2022. Cocoa-mt: A dataset and benchmark for contrastive controlled MT with application to formality. In Carpuat, Marine, Marie-Catherine de Marneffe, and Iván Vladimir Meza Ruíz, editors, *Findings of the Association for Computational Linguistics: NAACL 2022, Seattle, WA, United States, July 10-15, 2022*, pages 616–632. Association for Computational Linguistics.
- Niu, Xing and Marine Carpuat. 2020. Controlling neural machine translation formality with synthetic supervision. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 8568–8575. AAAI Press.
- Niu, Xing, Marianna J. Martindale, and Marine Carpuat. 2017. A study of style in machine translation: Controlling the formality of machine translation output. In Palmer, Martha, Rebecca Hwa, and Sebastian Riedel, editors, *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, pages 2814–2819. Association for Computational Linguistics.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Pavlick, Ellie and Joel Tetreault. 2016. An empirical analysis of formality in online communication. *Transactions of the Association for Computational Linguistics*, 4:61–74.
- Popović, Maja. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal, September. Association for Computational Linguistics.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October. Association for Computational Linguistics.
- Rao, Sudha and Joel R. Tetreault. 2018. Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. In Walker, Marilyn A., Heng Ji, and Amanda Stent, editors, *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 129–140. Association for Computational Linguistics.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.
- Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144.
- Schioppa, Andrea, David Vilar, Artem Sokolov, and Katja Filippova. 2021. Controlling machine translation for multiple attributes with additive interventions. In Moens, Marie-Francine, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 6676–6696. Association for Computational Linguistics.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2016. Controlling politeness in neural machine translation via side constraints. In Knight, Kevin, Ani Nenkova, and Owen Rambow, editors, *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 35–40. The Association for Computational Linguistics.
- Tiedemann, Jörg. 2016. OPUS – parallel corpora for everyone. In *Proceedings of the 19th Annual Conference of the European Association for Machine Translation: Projects/Products*, Riga, Latvia, May 30–June 1. Baltic Journal of Modern Computing.
- Tiedemann, Jörg. 2020. The Tatoeba Translation Challenge – Realistic data sets for low resource and multilingual MT. In *Proceedings of the Fifth Conference on Machine Translation*, pages 1174–1182, Online, November. Association for Computational Linguistics.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Guyon, I., U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and

R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Xue, Linting, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer. In Toutanova, Kristina, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 483–498. Association for Computational Linguistics.

Yang, Kevin and Dan Klein. 2021. FUDGE: controlled text generation with future discriminators. In Toutanova, Kristina, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 3511–3535. Association for Computational Linguistics.

Zaragoza-Bernabeu, Jaume, Gema Ramírez-Sánchez, Marta Bañón, and Sergio Ortiz Rojas. 2022. Bicleaner AI: Bicleaner goes neural. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 824–831, Marseille, France, June. European Language Resources Association.

A Detailed analysis of the FAME-MT dataset

To understand the **FAME-MT** dataset deeper, we include the extension of the explorative analysis by adding figures and tables that may help understand the dataset better.

Table 9 provides information about the size of each language pair subset in **FAME-MT**. For each language pair, 50% of examples (50,000) are examples with formal targets and 50% of examples with informal ones.

Table 10 extends the discussion on reading scores in terms of Flesch readability scores. Here, for each language pair, we see that informal translations are easier to read than formal ones.

Table 11 provides an in-depth analysis of the formality of source sentences in relation to the formality of target sentences extending information from Table 4. Here, we provide statistics for each language pair considering all kinds of disagreements separately, counting how many source sen-

tences were classified as formal/informal/neutral when target sentence is formal or informal.

Finally, Figures 2 and 3 provide information about the punctuation distribution and word length distribution for formal and informal texts. While formal documents tend to use longer words, the distribution of punctuation marks is quite similar for both categories.

Source language \ Target language	German (DE)	English (EN)	Spanish (ES)	French (FR)	Italian (IT)	Dutch (NL)	Polish (PL)	Portuguese (PT)
Czech (CS)	✓	✓	✓	✓	✓	✓	✓	✓
Danish (DA)	✓	✓	✓	✓	✓	✓	✓	✓
German (DE)	NOT	✓	✓	✓	✓	✓	✓	✓
English (EN)	✓	NOT	✓	✓	✓	✓	✓	✓
Spanish (ES)	✓	✓	NOT	✓	✓	✓	✓	✓
French (FR)	✓	✓	✓	NOT	✓	✓	✓	✓
Italian (IT)	✓	✓	✓	✓	NOT	✓	✓	✓
Dutch (NL)	✓	✓	✓	✓	✓	NOT	✓	✓
Norwegian (NO+NB)	✓	✓	✓	✓	✓	✓	✓	✓
Polish (PL)	✓	✓	✓	✓	✓	✓	NOT	✓
Portuguese (PT)	✓	✓	✓	✓	✓	✓	✓	NOT
Russian (RU)	✓	✓	✓	✓	✓	✓	✓	✓
Slovak (SK)	✓	✓	✓	✓	✓	✓	✓	✓
Swedish (SV)	✓	✓	✓	✓	✓	✓	✓	✓
Ukrainian (UK)	✓	✓	✓	✓	✓	✓	✓	✓

Table 9: Language pairs in FAME-MT. The total size of FAME-MT = 11.2 million examples (14 source languages · 8 target languages · 100,000 examples). Each cell with ✓ represents a language pair with 100,000 examples (50% formal translations and 50% informal). Each cell with NOT represents a language pair that is not supported because source language equals target language.

Source language \ Target language	German (DE)	English (EN)	Spanish (ES)	French (FR)	Italian (IT)	Dutch (NL)	Portuguese (PT)
Czech (CS)	53.88/ 77.84	63.55/ 83.98	75.05/ 88.91	82.41/ 103.41	78.95/ 88.3	60.71/ 84.91	75.93/ 92.6
Danish (DA)	50.76/ 72.02	64.37/ 86.56	65.67/ 78.72	76.9/ 100.84	66.95/ 80.14	50.15/ 73.15	62.53/ 89.26
German (DE)	-/-	61.56/ 83.35	66.11/ 76.03	76.54/ 95.52	65.74/ 75.74	56.96/ 72.34	61.56/ 85.69
English (EN)	49.81/ 63.57	-/-	63.49/ 72.15	76.55/ 95.33	67.05/ 75.14	57.84/ 72.29	63.55/ 83.46
Spanish (ES)	49.85/ 69.43	59.98/ 92.34	-/-	76.91/ 97.05	72.54/ 82.05	57.41/ 75.33	63.89/ 84.97
French (FR)	50.4/ 68.33	58.58/ 93.71	66.47/ 78.62	-/-	68.77/ 78.45	57.6/ 74.74	66.96/ 87.36
Italian (IT)	49.9/ 69.84	61.37/ 91.32	66.9/ 78.4	76.55/ 97.37	-/-	56.66/ 74.41	64.4/ 86.08
Dutch (NL)	51.31/ 71.75	67.11/ 94.6	69.67/ 81.12	77.03/ 101.03	70.41/ 82.1	-/-	68.43/ 90.56
Norwegian (NO+NB)	51.88/ 72.08	68.05/ 89.17	69.1/ 75.82	75.55/ 100.69	74.42/ 79.05	54.4/ 62.36	66.94/ 88.24
Polish (PL)	53.12/ 78.02	65.36/ 93.55	71.68/ 84.36	80.53/ 103.4	76.89/ 87.86	60.76/ 83.15	74.61/ 92.08
Portuguese (PT)	52.12/ 74.13	62.01/ 95.72	70.07/ 83.84	79.54/ 101.12	72.92/ 84.36	60.29/ 83.69	-/-
Russian (RU)	54.71/ 78.11	60.37/ 94.95	69.43/ 86.05	82.92/ 102.64	79.77/ 88.9	62.07/ 86.24	72.61/ 94.38
Slovak (SK)	52.15/ 73.61	62.26/ 90.41	66.31/ 81.35	77.27/ 99.28	65.89/ 78.68	56.78/ 76.96	63.21/ 88
Swedish (SV)	52.4/ 73.12	65.26/ 91.95	69.56/ 80.21	77.06/ 100.96	68.14/ 79.95	57.67/ 74.65	67.58/ 90.9
Ukrainian (UK)	55.63/ 74.16	64.76/ 90.88	70.1/ 80.79	80.79/ 96.98	74.69/ 78.82	57.27/ 74.29	64.38/ 84.51
Average	51.99/ 72.57	63.19/ 90.89	68.93/ 80.46	78.33/ 99.69	71.65/ 81.4	57.61/ 76.32	66.9/ 88.44

Table 10: Flesch readability scores between each pair of languages (higher value → text easier to read). Each cell contains two numbers separated by a slash sign. The first number represents the Flesch score calculated for a random sample of 5000 formal sentences expressed in a given language pair. The second number represents an analogous score for informal target sentences. In each sample, for each language pair, informal texts are easier to read (higher scores marked in bold).

Source and target languages	Formal target			Informal target		
	Formal source	Neutral source	Informal source	Informal source	Neutral source	Formal source
English (EN) → German (DE)	35070 (70.14%)	12390 (24.78%)	2540 (5.08%)	5281 (10.56%)	14562 (29.12%)	30157 (60.31%)
Spanish (ES) → German (DE)	17555 (35.11%)	21012 (42.02%)	11433 (22.87%)	32253 (64.51%)	11288 (22.58%)	6459 (12.92%)
French (FR) → German (DE)	33929 (67.86%)	15242 (30.48%)	829 (1.66%)	18925 (37.85%)	9963 (19.93%)	21112 (42.22%)
Italian (IT) → German (DE)	5496 (10.99%)	31903 (63.81%)	12601 (25.2%)	30364 (60.73%)	16867 (33.73%)	2769 (5.54%)
Polish (PL) → German (DE)	6498 (13.0%)	25771 (51.54%)	17731 (35.46%)	35236 (70.47%)	12164 (24.33%)	2600 (5.2%)
Portugal (PT) → German (DE)	29499 (59.0%)	19213 (38.43%)	1288 (2.58%)	12310 (24.62%)	12770 (25.54%)	24920 (49.84%)
Dutch (NL) → German (DE)	21803 (43.61%)	15463 (30.93%)	12734 (25.47%)	35478 (70.96%)	10576 (21.15%)	3946 (7.89%)
German (DE) → English (EN)	8444 (16.89%)	39829 (79.66%)	1727 (3.45%)	2783 (5.57%)	42584 (85.17%)	4633 (9.27%)
Spanish (ES) → English (EN)	8489 (16.98%)	36997 (73.99%)	4514 (9.03%)	9080 (18.16%)	35870 (71.74%)	5050 (10.1%)
French (FR) → English (EN)	8760 (17.52%)	40154 (80.31%)	1086 (2.17%)	5765 (11.53%)	36896 (73.79%)	7339 (14.68%)
Italian (IT) → English (EN)	4204 (8.41%)	40544 (81.09%)	5252 (10.5%)	7748 (15.5%)	39355 (78.71%)	2897 (5.79%)
Polish (PL) → English (EN)	2821 (5.64%)	38964 (77.93%)	8215 (16.43%)	14810 (29.62%)	32751 (65.5%)	2439 (4.88%)
Portugal (PT) → English (EN)	13104 (26.21%)	35920 (71.84%)	976 (1.95%)	3234 (6.47%)	35328 (70.66%)	11438 (22.88%)
Dutch (NL) → English (EN)	5585 (11.17%)	37794 (75.59%)	6621 (13.24%)	8831 (17.66%)	39525 (79.05%)	1644 (3.29%)
German (DE) → Spanish (ES)	21142 (42.28%)	26750 (53.5%)	2108 (4.22%)	17301 (34.6%)	10018 (20.04%)	22681 (45.36%)
English (EN) → Spanish (ES)	36314 (72.63%)	10132 (20.26%)	3554 (7.11%)	8148 (16.3%)	14101 (28.2%)	27751 (55.5%)
French (FR) → Spanish (ES)	16087 (32.17%)	32477 (64.95%)	1436 (2.87%)	13417 (26.83%)	10920 (21.84%)	25663 (51.33%)
Italian (IT) → Spanish (ES)	15430 (30.86%)	27140 (54.28%)	7430 (14.86%)	27970 (55.94%)	18995 (37.99%)	3035 (6.07%)
Polish (PL) → Spanish (ES)	7099 (14.2%)	30481 (60.96%)	12420 (24.84%)	34870 (69.74%)	12181 (24.36%)	2949 (5.9%)
Portugal (PT) → Spanish (ES)	34107 (68.21%)	14891 (29.78%)	1002 (2.0%)	9874 (19.75%)	12240 (24.48%)	27886 (55.77%)
Dutch (NL) → Spanish (ES)	14564 (29.13%)	27806 (55.61%)	7630 (15.26%)	31427 (62.85%)	9942 (19.88%)	8631 (17.26%)
German (DE) → French (FR)	36371 (72.74%)	7966 (15.93%)	5663 (11.33%)	29805 (59.61%)	14752 (29.5%)	5443 (10.89%)
English (EN) → French (FR)	33675 (67.35%)	12696 (25.39%)	3629 (7.26%)	15148 (30.3%)	12471 (24.94%)	22381 (44.76%)
Spanish (ES) → French (FR)	17940 (35.88%)	15202 (30.4%)	16858 (33.72%)	30891 (61.78%)	13582 (27.16%)	5527 (11.05%)
Italian (IT) → French (FR)	5948 (11.9%)	27005 (54.01%)	17047 (34.09%)	29517 (59.03%)	16683 (33.37%)	3800 (7.6%)
Polish (PL) → French (FR)	8272 (16.54%)	18179 (36.36%)	23549 (47.1%)	36702 (73.4%)	14101 (28.2%)	2421 (4.84%)
Portugal (PT) → French (FR)	33897 (67.79%)	13684 (27.37%)	2419 (4.84%)	14252 (28.5%)	14079 (28.16%)	21669 (43.34%)
Dutch (NL) → French (FR)	23659 (47.32%)	7537 (15.07%)	18804 (37.61%)	34132 (68.26%)	14393 (28.79%)	1475 (2.95%)
German (DE) → Italian (IT)	17323 (34.65%)	31410 (62.82%)	1267 (2.5%)	15333 (30.67%)	8605 (17.21%)	26062 (52.12%)
English (EN) → Italian (IT)	34491 (68.98%)	10309 (20.62%)	5200 (10.4%)	8882 (17.76%)	12870 (25.74%)	28248 (56.5%)
Spanish (ES) → Italian (IT)	28522 (57.04%)	16438 (32.88%)	5040 (10.08%)	30705 (61.41%)	9555 (19.11%)	9740 (19.48%)
French (FR) → Italian (IT)	11666 (23.33%)	36215 (72.43%)	2119 (4.24%)	13321 (26.64%)	10083 (20.17%)	26596 (53.19%)
Polish (PL) → Italian (IT)	9932 (19.86%)	27943 (55.89%)	12125 (24.25%)	34979 (69.96%)	12180 (24.36%)	2841 (5.68%)
Portugal (PT) → Italian (IT)	31358 (62.72%)	17083 (34.17%)	1559 (3.12%)	10006 (20.01%)	11631 (23.26%)	28363 (56.73%)
Dutch (NL) → Italian (IT)	8003 (16.01%)	34786 (69.57%)	7211 (14.42%)	29321 (58.64%)	9900 (19.8%)	10779 (21.56%)
German (DE) → Polish (PL)	20330 (40.66%)	26744 (53.49%)	2926 (5.85%)	12171 (24.34%)	20941 (41.88%)	16888 (33.78%)
English (EN) → Polish (PL)	31044 (62.09%)	10617 (21.23%)	8339 (16.68%)	13598 (27.2%)	11868 (23.74%)	24534 (49.07%)
Spanish (ES) → Polish (PL)	16357 (32.71%)	28082 (56.16%)	5561 (11.12%)	18445 (36.89%)	18661 (37.32%)	12894 (25.79%)
French (FR) → Polish (PL)	19355 (38.71%)	27293 (54.59%)	3352 (6.7%)	13362 (26.72%)	21551 (43.1%)	15087 (30.17%)
Italian (IT) → Polish (PL)	11410 (22.82%)	32094 (64.19%)	6496 (12.99%)	20093 (40.19%)	26253 (52.51%)	3654 (7.31%)
Portugal (PT) → Polish (PL)	20408 (40.82%)	27196 (54.39%)	2396 (4.79%)	8348 (16.7%)	22168 (44.34%)	19484 (38.97%)
Dutch (NL) → Polish (PL)	11248 (22.5%)	28816 (57.63%)	9936 (19.87%)	22401 (44.8%)	22483 (44.97%)	5116 (10.23%)
German (DE) → Portugal (PT)	22322 (44.64%)	20169 (40.34%)	7509 (15.02%)	27060 (54.12%)	15757 (31.51%)	7183 (14.37%)
English (EN) → Portugal (PT)	33023 (66.05%)	11556 (23.11%)	5421 (10.84%)	14020 (28.04%)	12330 (24.66%)	23650 (47.3%)
Spanish (ES) → Portugal (PT)	21487 (42.97%)	14667 (29.33%)	13846 (27.69%)	31533 (63.07%)	14485 (28.97%)	3982 (7.96%)
French (FR) → Portugal (PT)	21263 (42.53%)	22366 (44.73%)	6371 (12.74%)	25245 (50.49%)	15808 (31.62%)	8947 (17.89%)
Italian (IT) → Portugal (PT)	10367 (20.73%)	24452 (48.9%)	15181 (30.36%)	29337 (58.67%)	17779 (35.56%)	2884 (5.77%)
Polish (PL) → Portugal (PT)	5696 (11.39%)	20486 (40.97%)	23818 (47.64%)	35852 (71.7%)	11792 (23.58%)	2356 (4.71%)
Dutch (NL) → Portugal (PT)	10915 (21.83%)	19712 (39.42%)	19373 (38.75%)	33186 (66.37%)	15465 (30.93%)	1349 (2.7%)
German (DE) → Dutch (NL)	40742 (81.48%)	7124 (14.25%)	2134 (4.27%)	16806 (33.61%)	11831 (23.66%)	21363 (42.73%)
English (EN) → Dutch (NL)	36558 (73.12%)	10737 (21.47%)	2705 (5.41%)	8698 (17.4%)	12496 (24.99%)	28806 (57.61%)
Spanish (ES) → Dutch (NL)	22760 (45.52%)	15556 (31.11%)	11684 (23.37%)	28258 (56.52%)	13826 (27.65%)	7916 (15.83%)
French (FR) → Dutch (NL)	41421 (82.84%)	7746 (15.49%)	833 (1.67%)	14766 (29.53%)	9928 (19.86%)	25306 (50.61%)
Italian (IT) → Dutch (NL)	6583 (13.17%)	28871 (57.74%)	14546 (29.09%)	25372 (50.74%)	20634 (41.27%)	3994 (7.99%)
Polish (PL) → Dutch (NL)	10729 (21.46%)	21618 (43.24%)	17653 (35.31%)	33024 (66.05%)	13191 (26.38%)	3785 (7.57%)
Portugal (PT) → Dutch (NL)	35708 (71.42%)	13140 (26.28%)	1152 (2.3%)	11072 (22.14%)	13441 (26.88%)	25487 (50.97%)

Table 11: Relations between source sentence formality and target sentence formality in FAME-MT determined using our classifiers. For each formal and informal target sentence, the classifier is used to determine the formality of the corresponding source sentence. Then, the number of source sentences classified as formal, informal, and neutral is reported for each target category.

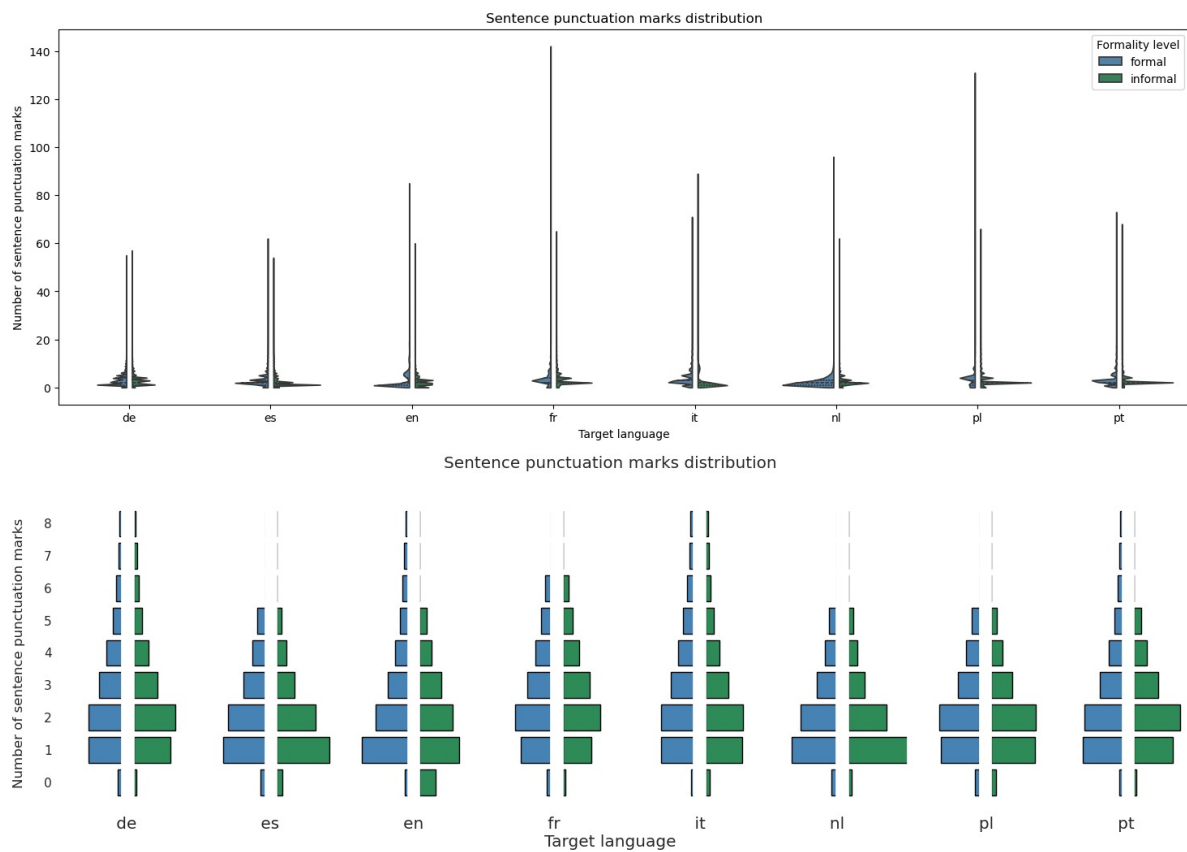


Figure 2: Plots representing the distributions of the number of punctuation signs in a sentence for a given language. The upper figure represents the distributions calculated over the original dataset. As it shows that there are some outliers, we provide the lower figure generated over a subset of texts whose lengths are between $Q1 - 1.5 IQR$ and $Q3 + 1.5 IQR$ ($Q1$ =first quartile, $Q3$ =third quartile, IQR =inter-quartile range) to focus more on the most common scenarios.

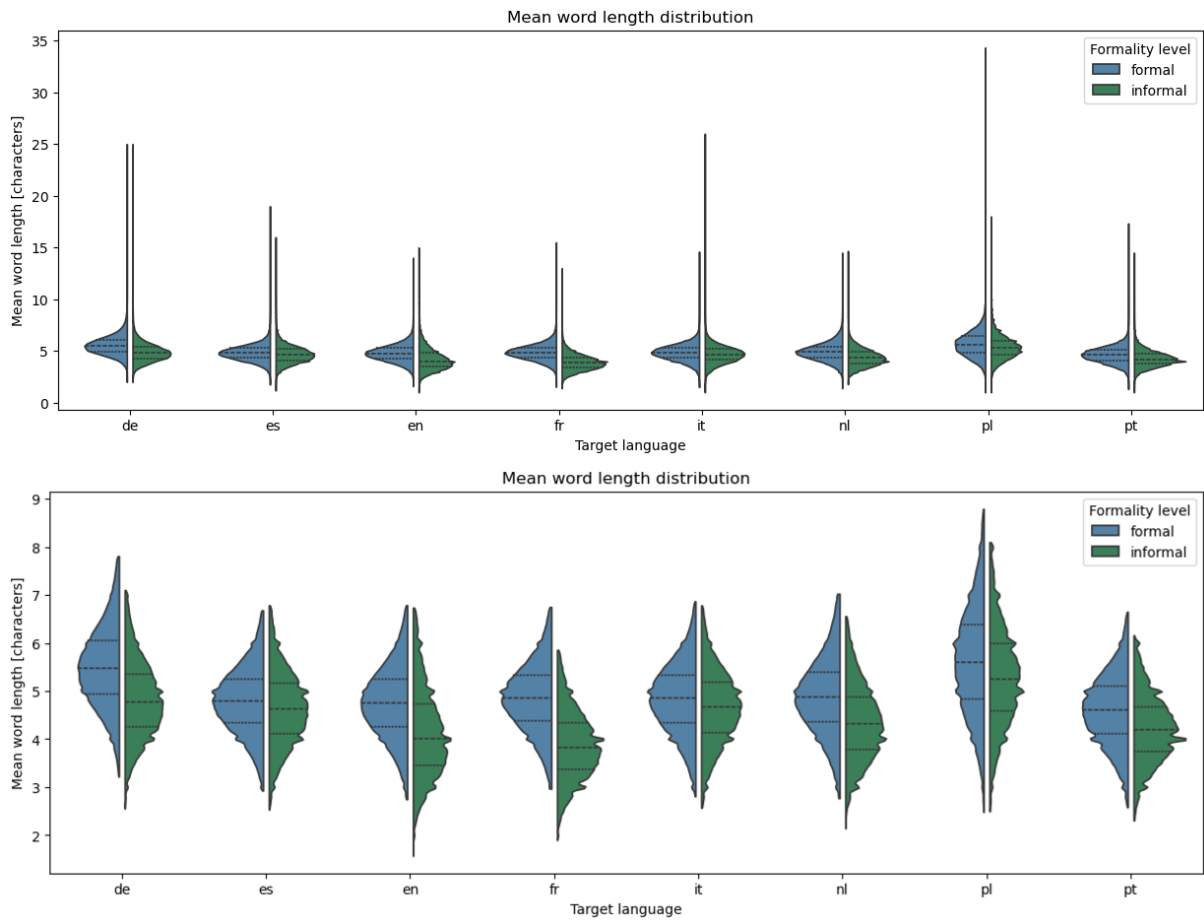


Figure 3: Plots representing the distributions of the mean word length in a given sentence per given language. The upper figure represents the distributions calculated over the original dataset. As it shows that there are some outliers, we provide the lower figure generated over a subset of texts whose lengths are between $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$ ($Q1$ =first quartile, $Q3$ =third quartile, IQR =inter-quartile range) to focus more on the most common scenarios.

Iterative Translation Refinement with Large Language Models

Pinzhen Chen¹ Zhicheng Guo² Barry Haddow¹ Kenneth Heafield¹

¹School of Informatics, University of Edinburgh

²Dept. of Comp. Sci. & Tech., Institute for AI, Tsinghua University

{pinzhen.chen, bhaddow, kenneth.heafield}@ed.ac.uk

guo-zc21@mails.tsinghua.edu.cn

Abstract

We propose iteratively prompting a large language model to self-correct a translation, with inspiration from their strong language capability as well as a human-like translation approach. Interestingly, multi-turn querying reduces the output’s string-based metric scores, but neural metrics suggest comparable or improved quality after two or more iterations. Human evaluations indicate better fluency and naturalness compared to initial translations and even human references, all while maintaining quality. Ablation studies underscore the importance of anchoring the refinement to the source and a reasonable seed translation for quality considerations. We also discuss the challenges in evaluation and relation to human performance and translationese.

1 Introduction

Large language models (LLMs), e.g. generative pre-trained Transformers (GPT), have made notable advancements in natural language processing (Radford et al., 2019; Brown et al., 2020; Kaplan et al., 2020; Ouyang et al., 2022). In machine translation (MT), where the convention is to use an encoder-decoder architecture to deal with source and target sentences respectively (Bahdanau et al., 2015; Vaswani et al., 2017), recent papers have examined the feasibility of LLM prompting for translation (Vilar et al., 2023; Zhang et al., 2023; Hendy et al., 2023; Agrawal et al., 2023).

With autoregressive decoding being the convention, machine translation models yield output in

a single attempt, and so do post-editing models. Rather, a human translator can read and edit translations repeatedly, or even pass the outcome to another translator for a second opinion. We explore such an iterative refinement process with LLMs, where the proposed method simply feeds a source-translation pair into an LLM for an improved translation in multiple rounds. It is worth noting that this method can be applied to an initial translation from any model, not just LLM outputs. We further conduct a qualitative evaluation of the outputs. Our approach offers two insights from a fluency and naturalness perspective: 1) LLMs are pre-trained on natural texts that are orders of magnitude larger than traditional MT data, and 2) the method does not require complicated prompt engineering, yet allows for iterative and arbitrary rephrasing compared to automatic post-editing, which is limited to token-level error correction without style editing (Ive et al., 2020).

Empirical results show that the refinement procedure introduces significant textual changes reflected by the drop in BLEU and chrF++, but attains similar or higher COMET scores compared to initial translations. Native speakers prefer refined outputs in terms of fluency and naturalness when compared with GPT translations and even human references. Reference-based human evaluation confirms that such gains are made without sacrificing general quality. As corroborated by recent works, automatic metrics like BLEU and COMET are witnessed to move in opposite directions (Freitag et al., 2019; Freitag et al., 2022). Our human-like LLM prompting method contributes to translation naturalness which can enhance utility as perceived by the target language users. On a broader scope, this work touches on the concept of involving LLMs in a collaborative translation editing strategy.

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

Mode	Prompt
<i>Translate</i>	Source: $\{\text{source}\}$ Please give me a translation in $\{\text{lang}\}$ without any explanation.
<i>Refine</i>	Source: $\{\text{source}\}$ Translation: $\{\text{prev_translation}\}$ Please give me a better $\{\text{lang}\}$ translation without any explanation.
<i>Refine_{Contrast}</i>	Source: $\{\text{source}\}$ Bad translation: $\{\text{prev_translation}\}$ Please give me a better $\{\text{lang}\}$ translation without any explanation.
<i>Refine_{Random}</i>	Source: $\{\text{source}\}$ Bad translation: $\{\text{random_target}\}$ if first-round, else $\{\text{prev_translation}\}$ Please give me a better $\{\text{lang}\}$ translation without any explanation.
<i>Paraphrase</i>	Sentence: $\{\text{prev_translation}\}$ Please give me a paraphrase in $\{\text{lang}\}$ without any explanation.

Table 1: Prompts used in our work, where a $\{\text{variable}\}$ is substituted with its corresponding content.

2 Methodology

Having an input source sentence x and an optimizable model θ_{mt} , the process to obtain a translation y can be modelled as $y = \text{argmax}_y P(y|x; \theta_{mt})$. Next, an automatic post-editor θ_{ape} creates a refined translation y' through modelling $y' = \text{argmax}_{y'} P(y'|x, y; \theta_{ape})$. Conventional translation or automatic post-editing models are trained on (x, y) or (x, y, y') data pairs.

Extending prior work on LLM prompting, our study uses zero-shot prompting by affixing a task description to form a prompt p and querying an LLM θ_{LLM} to elicit a response (Brown et al., 2020). We introduce five prompts in our study:

1. *Translate*: it queries for a translation of a source input, extending the translation process with a prompt p : $y = \text{argmax}_y P(y|p, x; \theta_{LLM})$. This is vanilla LLM prompting for MT.
2. *Refine*: similar to post-editing, the LLM is given the source sentence and the previous translation to produce a better translation $y' = \text{argmax}_{y'} P(y'|p, x, y; \theta_{LLM})$.
3. *Refine_{Contrast}*: as a contrasting prompt to the above, we insert the word “bad” to hint that the previously translated text is unwanted, regardless of its actual quality.
4. *Refine_{Random}*: same prompt as *Refine_{Contrast}*, but in the first iteration, a random sentence is fed instead of a translation to imitate a genuinely “bad translation”.
5. *Paraphrase*: a contrasting experiment to translation prompting, we ask an LLM to rephrase a translation without feeding the source sentence x : $y'' = \text{argmax}_{y''} P(y''|p, y; \theta_{LLM})$.

We propose to iteratively call the refinement prompts, where the source stays the same but the previous translation is updated each turn. To encourage a parsable model response, we ask the LLM to not give any explanation. Such prompting does not require model parameters θ_{LLM} to be accessible. Through ablation prompts, *Refine_{Random}* and *Paraphrase*, we analyse to what degree the source input and seed translations are helpful. The exact prompt texts are displayed in Table 1.

3 Experiments

3.1 Data and model details

We select language pairs from the news and general domain translation tasks hosted at WMT 2021 and 2022 (Farhad et al., 2021; Kocmi et al., 2022), which are supported by COMET to obtain reliable scores. In total, we tested seven translation directions: English \leftrightarrow German (en \rightarrow de, de \rightarrow en), English \leftrightarrow Chinese (en \rightarrow zh, zh \rightarrow en), German \rightarrow French (de \rightarrow fr), English \rightarrow Japanese (en \rightarrow ja), and Ukrainian \rightarrow Czech (uk \rightarrow cs). We directly benchmark on the test sets, and in situations where multiple references are available, we use human reference “A” released by the WMT organizers as our reference.

We experiment with GPT-3.5, a powerful closed-source model from OpenAI that can be accessed by all users.¹ As the API call tends to be slow, we randomly sample 200 instances from the official test set to form our in-house test. In the refinement and paraphrase experiments, we use the response from

¹We accessed a version of gpt-3.5-turbo with training data up to Sep 2021, so it should not have seen WMT 2021 or 2022 test references. Nevertheless, our findings are mostly drawn from reference-free metrics and human evaluation.

the LLM *Translate* query as the seed translation to be improved upon. We do not keep the query (multi-turn) history so as to prevent an LLM from seeing that the previous translation is produced by itself. In experiments later on, we also tested with translations from encoder-decoder systems that participated in WMT, human references, and online systems. Overall, translation refinement is iterated four times at maximum considering the API costs.

3.2 Evaluation setup

We consider four automatic metrics: string-based BLEU (Papineni et al., 2002) and chrF++ (Popović, 2017) as well as embedding-based COMET_{DA} and COMET_{QE} (Rei et al., 2020). The difference between the DA and QE versions is that COMET_{DA} requires a source, a translation, and a reference, whereas COMET_{QE} is reference-free. BLEU and chrF++ are as implemented in the sacrebleu toolkit.² We also use this toolkit to obtain test sets with references as well as past WMT systems’ outputs. Specifically for tokenization in BLEU calculation, we use “zh” for Chinese, “ja-mecab” for Japanese, and “13a” for the rest. The BLEU and chrF++ signatures are footnoted.^{3,4} For COMET metrics, we used the official implementation released by the authors.⁵

3.3 Refinement results

WMT21 We first experiment with en↔de and en↔zh from WMT21, which are high-resource languages in terms of both translation data and LLM training data. We run all five prompts and display results in Table 2. For iterative refinement and paraphrasing experiments, the best iteration is picked according to COMET_{QE}. We observe that the refined translations record a drastic drop in string-based metrics compared to initial translations, indicating lexical and structural variations. In terms of COMET_{DA}, refined outputs surpass initial GPT translations in three out of four cases, and in terms of COMET_{QE}, the refinement strategy ends as the highest with substantial improvement for into-English directions. As a contrasting experiment, *Paraphrase* sees a decline in all metrics, suggesting the importance of feeding the source input as an anchor during iterations to prevent semantic drift.

²<https://github.com/mjpost/sacrebleu>

³#:1|c:mixed|e:no|tok:13a|s:exp|v:2.3.1

⁴#:1|c:mixed|e:yes|nc:6|nw:2|s:no|v:2.3.1

⁵<https://github.com/Unbabel/COMET>

	BLEU	chrF++	COMET _{DA}	COMET _{QE}
<i>Reference_A</i>	-	-	-	.0919
de Translate	30.90	57.55	.8606	.1128
↓ Refine	23.14	51.91	.8525	.1116
en Refine _{Contrast}	22.88	52.47	.8452	.1162
Refine _{Random}	18.83	51.79	.7777	.0770
Paraphrase	11.01	40.05	.8044	.0919
<i>Reference_A</i>	-	-	-	.1127
en Translate	25.39	53.54	.8427	.1083
↓ Refine	22.35	50.57	.8478	.1153
de Refine _{Contrast}	22.54	51.21	.8211	.0929
Refine _{Random}	19.36	46.56	.7906	.0832
Paraphrase	13.60	43.54	.8197	.1006
<i>Reference_A</i>	-	-	-	.0708
zh Translate	25.64	53.74	.8199	.0867
↓ Refine	20.26	49.06	.8156	.0921
en Refine _{Contrast}	24.81	51.77	.8538	.1132
Refine _{Random}	24.24	47.11	.8323	.1022
Paraphrase	12.76	40.92	.7931	.0885
<i>Reference_A</i>	-	-	-	.0956
en Translate	29.28	20.61	.8300	.0761
↓ Refine	28.26	19.28	.8417	.0870
zh Refine _{Contrast}	29.28	19.69	.8395	.0881
Refine _{Random}	25.71	17.49	.8126	.0763
Paraphrase	21.95	17.14	.8144	.0716

Table 2: Automatic scores of different strategies with GPT on high-resource pairs from WMT 2021 news translation.

WMT22 Moving to lower-resourced languages with non-English translation, we gather numbers for three translation directions from WMT22 in Table 3. Since *Refine_{Random}* results are not desirable for WMT21, we omit experiments with this. The overall pattern remains the same as before: *Refine* works best, obtaining higher COMET_{QE} than vanilla translations and *Refine_{Contrast}*. Also, the reduction in string-based scores becomes less obvious, which might be attributed to seed GPT translations in lesser-resourced languages being lower in quality in the beginning.

Online systems, encoder-decoder systems, and human translations

In addition to translation refinement from GPT-3.5 itself, we also apply our refinement calls to outputs from conventional MT systems and human translators. These translations can represent genuine errors, if any, introduced during the translation process. Out of the seven WMT21 submissions, we select outputs from four models built by research labs that, based on human evaluation, have been ranked at significantly different positions on the German-to-English leaderboard: Tencent (Wang et al., 2021), Facebook AI (Tran et al., 2021), Edinburgh (Chen et al., 2021),

		BLEU	chrF++	COMET _{DA}	COMET _{QE}
	<i>Reference</i>	-	-	-	.0772
de ↓ fr	Translate	36.25	59.50	.8395	.0807
	Refine	32.47	55.83	.8353	.0851
	Refine _{Contrast}	33.12	56.37	.8308	.0805
	Paraphrase	16.06	44.28	.7937	.0682
	<i>Reference</i>	-	-	-	.1345
en ↓ ja	Translate	23.00	25.89	.8863	.1255
	Refine	22.63	27.30	.8941	.1305
	Refine _{Contrast}	22.82	26.71	.8928	.1282
	Paraphrase	17.69	23.18	.8592	.1086
	<i>Reference</i>	-	-	-	.1273
uk ↓ cs	Translate	29.91	54.64	.9074	.1173
	Refine	28.60	53.06	.9040	.1183
	Refine _{Contrast}	28.90	54.29	.9036	.1151
	Paraphrase	13.59	40.04	.8625	.0969

Table 3: Automatic scores of different strategies with GPT on low-resource and medium-resource pairs from WMT 2022 news translation.

and Huawei TSC (Wei et al., 2021). These are competitive systems built with data augmentation, multilingualism, ensembling, re-ranking, etc. We then include two online engines used in WMT 2021: Online-A and Online-Y. Finally, human reference “B” is added so that we can experiment with our refinement strategy on human translations.⁶ References “A” and “B” are sourced from different translation agencies (Farhad et al., 2021).

We report automatic scores from the refinement process in Table 4. A pattern similar to previous GPT translation refinement is noticed: for five out of seven WMT entries, the refinement strategy reaches a higher COMET_{QE} score, surprisingly, with up to one-third drop in BLEU. *Refine_{Contrast}* in all but one system surpass *Refine*, and without the initial translation, *Paraphrase* iterations record the lowest scores compared to the original submissions and refinements.

4 Human Evaluation

String-based and neural scores are observed to vary in opposite directions, which may suggest volatile changes in texts. Since it is questionable to conclude a quality degradation in this case, we set up human evaluations to measure two characteristics in the refined translations: text naturalness and overall quality. Human evaluators involved in this study

⁶The overview paper of WMT 2021 states that “for German↔English, the ‘B’ reference was found to be a post-edited version of one of the participating online systems”. We discover that it refers to English→German only, and German→English is not affected.

		BLEU	chrF++	COMET _{DA}	COMET _{QE}
	<i>Reference_A</i>	-	-	-	.0919
de ↓ fr	Submission	30.05	56.00	.8497	.1050
	Refine	23.39	51.80	.8527	.1123
	Refine _{Contrast}	25.10	53.82	.8566	.1116
	Paraphrase	12.52	41.03	.8031	.0894
	<i>Reference_B</i>	-	-	-	-
Online _A	Submission	34.45	60.78	.8582	.1061
	Refine	23.37	51.67	.8494	.1098
	Refine _{Contrast}	25.14	52.84	.8534	.1137
	Paraphrase	12.22	41.34	.8097	.0942
	<i>Reference_C</i>	-	-	-	-
Online _Y	Submission	32.70	59.32	.8500	.0981
	Refine	22.92	50.85	.8522	.1080
	Refine _{Contrast}	24.40	53.32	.8517	.1134
	Paraphrase	11.97	40.29	.8054	.0892
	<i>Reference_D</i>	-	-	-	-
Tencent	Submission	35.35	61.28	.8584	.1055
	Refine	23.75	52.16	.8488	.1095
	Refine _{Contrast}	26.89	54.75	.8553	.1116
	Paraphrase	12.43	41.35	.8116	.0947
	<i>Reference_E</i>	-	-	-	-
Facebook	Submission	34.67	60.78	.8677	.1146
	Refine	22.97	51.05	.8505	.1113
	Refine _{Contrast}	25.74	53.88	.8548	.1130
	Paraphrase	11.80	40.99	.8099	.0922
	<i>Reference_F</i>	-	-	-	-
Edinburgh	Submission	34.20	60.03	.8588	.1087
	Refine	22.04	50.29	.8496	.1097
	Refine _{Contrast}	25.24	52.87	.8546	.1147
	Paraphrase	12.79	40.18	.8067	.0921
	<i>Reference_G</i>	-	-	-	-
Huawei	Submission	35.13	61.17	.8643	.1126
	Refine	22.24	50.82	.8519	.1097
	Refine _{Contrast}	24.95	52.47	.8560	.1124
	Paraphrase	12.20	40.74	.8078	.0909

Table 4: Automatic scores of refining WMT 2021 news shared task German-to-English submissions.

are practitioners in the field of natural language processing but are unaware of the goal of this study.

4.1 Fluency and naturalness

We mimic the human evaluation of fluency in (Lembersky et al., 2012, p819). Native speakers of the target language are with two translations but without the source sentence; then we ask “Please choose the translation that is more fluent, natural, and reflecting better use of $\{\text{language}\}$ ”, where $\{\text{language}\}$ is substituted with the target language name. The evaluator has three options: they can select one of the two translations, or a “tie” if they consider both equally (un)natural. We conduct such pairwise evaluation to compare the first-round output from *Refine_{Contrast}* against human references, as well as against *Translate* separately.

We evaluate 50 samples from en↔de and en↔zh experiments in Section 3.3, and report in Figure 1 (left). Native speakers prefer *Refine_{Contrast}* to vanilla *Translate* in all four directions, and even favour

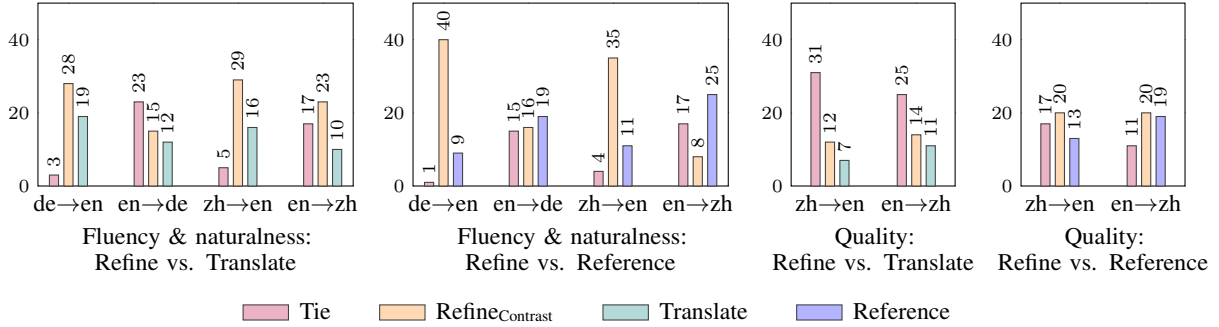


Figure 1: Human preferences on fluency and naturalness (source-free, left) and overall quality (source-based, right).

*Refine*_{Contrast} over human references when translating into English. It demonstrates that our simple strategy enhances the naturalness of GPT outputs and that WMT human references could be less favourable than GPT outputs in some cases.

4.2 Overall quality

We also evaluate for general quality as a safeguard. In this setup, a source sentence and two translations are given to an evaluator who is fluent in both languages. They are asked to pick the translation with better quality or indicate a tie. We only evaluated two translation directions, English to and from Chinese, due to the limited availability of bilingual speakers. Similar to the previous evaluation, we compare *Refine*_{Contrast} against human references, as well as *Refine*_{Contrast} against *Translate* separately.

We report evaluator preferences in Figure 1 (right). It shows that GPT *Refine* attains slightly better performance in zh→en and similar performance in en→zh when compared with human references. On the other hand, it is more favourable than GPT *Translate* in terms of human judgements. Combining evaluation outcomes, we conclude that the refinement strategy could improve the target-side naturalness without undermining general quality.

5 Analysis and Discussions

5.1 Performance through iterations

To investigate the behaviour of refinement strategies through different iterations, we plot BLEU, COMET_{DA}, and COMET_{QE} at different iterations in Figure 2 for four translation directions: en↔de and en↔zh. We find that *Refine* and *Refine*_{Contrast} usually attain their best after undergoing more than one refinement iteration, showing superiority to one-off editing.⁷ However, in almost all *Paraphrase*

⁷The first iteration is equivalent to a one-off translation editing using an LLM.

experiments, scores decrease monotonically, indicating that semantics drift away as paraphrasing iterates. Moreover, *Refine*_{Random} results start low, gradually catch up, but never reach as high as *Refine* or *Refine*_{Contrast}. This means that iterative refinement is indeed useful in fixing translations, but starting with a reasonable translation is also crucial for obtaining a strong result.

5.2 Diverging automatic scores

According to automatic string-based metrics, our queries deliver lower-quality translations through iterations, but COMET_{DA} scores remain comparable and COMET_{QE} scores mostly increase. We argue that the string-based metrics might not accurately indicate quality, but rather reflect text variations with respect to the reference. We further verified this via human evaluation that fluency and overall quality are not impacted.

In Table 5 we show outputs from different strategies for a single source input, where a native speaker marked preference for *Refine*_{Contrast}. It illustrates that the word choice is diverse for both directions and specifically for Chinese→English, there are substantial structural changes. The huge variety in expressions across translations can result in low BLEU with respect to human references, but without much change in meaning, for instance, as in Table 2 where BLEU can decline up to one-third, but neural metric scores change little. In the field of MT, a leap in BLEU is usually associated with performance improvement; however, in our case, a drop cannot be simply interpreted as performance degradation. This can be attributed to the lexical and structural diversity in the refined translations.

5.3 Human performance

A human translator is deemed to be fluent in their native language, which intuitively is difficult for a model to compete with. In our human evalua-

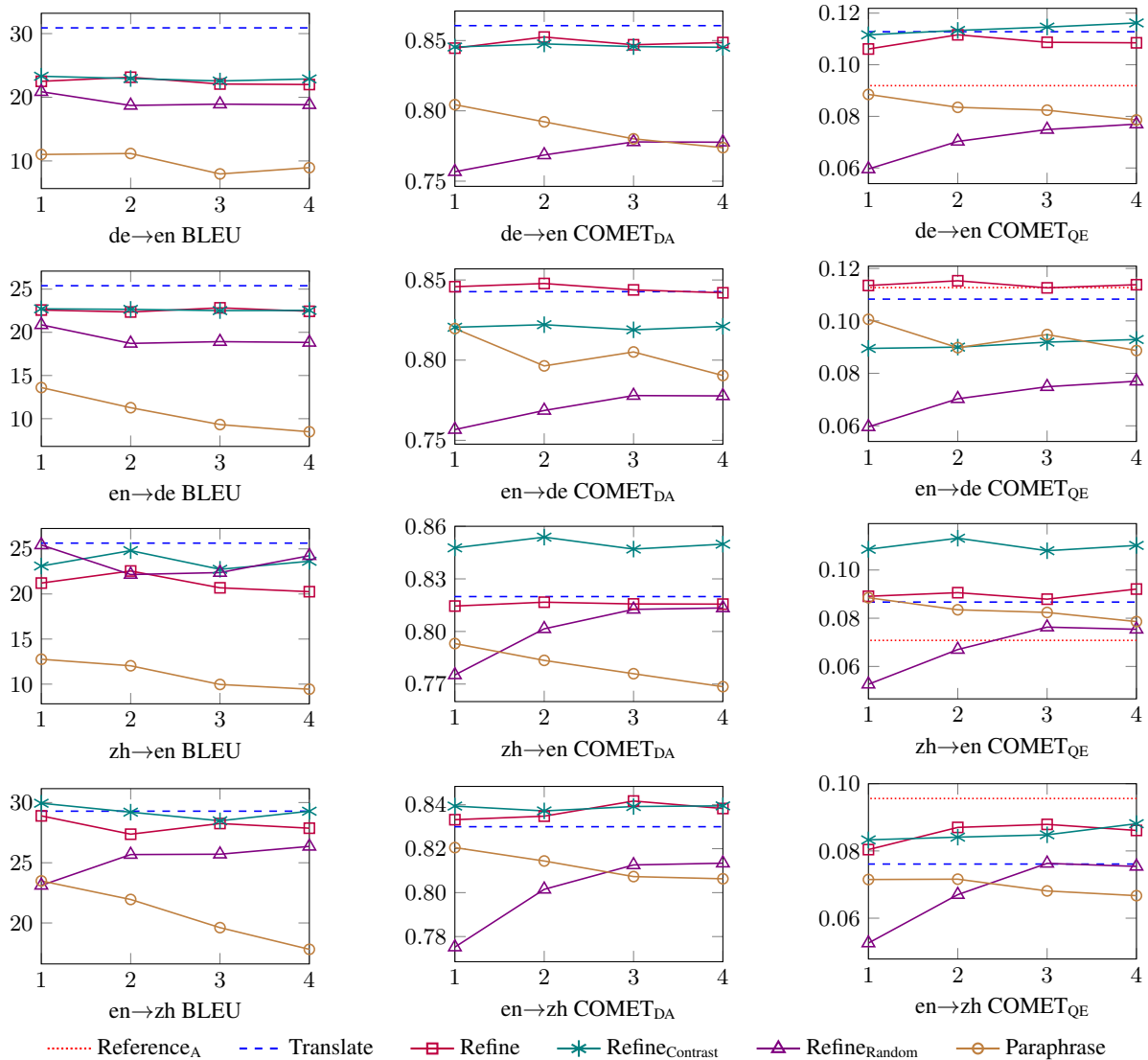


Figure 2: BLEU, COMET_{DA}, and COMET_{QE} at different refinement and paraphrase iterations for high-resource translation.

tion, GPT fluency can be as good or even better than reference translations—we offer two possible explanations. First, the WMT references might have been created by translators with varying expertise, which may not represent upper-bound human performance, especially when compared with advanced LLMs. More importantly, translations can exhibit awkwardness in word and syntax choices, potentially due to source language interference or “shining through” (Gellerstam, 1986; Teich, 2003).

5.4 Relation to translationese

Both human and machine translations might be more explicit, language-normalized, and simpler (Baker, 1996; Koppel and Ordan, 2011). On a broader scope, translationese is regarded as the distinct features in translations to include influences from both the source and target sides. Although

MT normally learns from human translation data, researchers found that human and machine translation patterns do not fully overlap (Bizzoni et al., 2020). While translationese occurs in translations inevitably, consumers could prefer translations that are more natural in their native language, provided that the semantics and utility are preserved.

From a narrow aspect, our method relates to machine translationese mitigation in terms of reducing unnaturalness and literalness, instead of focusing on state-of-the-art metric scores. It may be viable to create diverse translations through iterations, as we observe huge changes in BLEU scores. Measuring these using automatic metrics at the moment is challenging, especially given that most translation metrics are reference-based, where the reference can be translationese-prone in the first place. COMET_{QE} might be more robust to this end.

Source	Der 17-Jährige floh zunächst vom Tatort , seine Personalien konnten aber im Nachhinein ermittelt werden .
Reference	The 17 year-old proceeded to flee the crime scene , however, his personal details could be retrieved later.
Translate	The 17-year-old initially fled from the crime scene , but his personal information was later determined .
Refine _{Contrast}	The 17-year-old initially fled from the scene of the crime , but his personal details could later be identified .
Paraphrase	At first , the 17-year-old ran away from where the crime occurred , but eventually, the authorities were able to identify him by his personal details.
Source	新法令规定，坎帕尼亚大区自即日起室内公共场所必须戴口罩， 违者 最高可处以1000欧元罚金。
Reference	According to a new decree , people must wear masks in indoor public places in Campania from now on, and offenders can be fined up to 1,000 euros.
Translate	A new regulation stipulates that in Campania, indoor public places must wear masks. Violators can be fined up to 1000 euros.
Refine _{Contrast}	A new regulation states that in the Campania region, masks must be worn in indoor public places, with a maximum fine of 1000 euros for those who violate the rule .
Paraphrase	A new rule in Campania requires people to wear masks in indoor public places, and those who don't follow this rule may be charged up to 1000 euros.

Table 5: German→English and Chinese→English examples showing rich lexical variations across translation strategies.

6 Related Work

6.1 Translation post-editing

Closely related to our refinement prompting is automatic post-editing (APE), which trains a neural network to fix translation errors by learning from human correction data, that can be traced back to as early as (Knight and Chander, 1994). While it has shown advancements in statistical machine translation, it has been suspected to be less effective in the deep learning era due to original translations being high-quality and lack of post-editing data (Junczys-Dowmunt and Grundkiewicz, 2018; Chatterjee et al., 2018). Whilst one way to facilitate this is more data provision (Chollampatt et al., 2020; Ive et al., 2020), our workaround utilizes a large language model, which possesses the post-editing capability without the need for specific training or fine-tuning. Furthermore, post-editing models might have limited power to alleviate awkwardness, because human editing data is collected from annotators who are usually instructed to not make style improvements (Ive et al., 2020). Compared to APE, our method allows LLMs to re-generate an entirely different translation, which could escape the “post-editese” phenomenon, where Toral (2019) demonstrated that human-edited machine translations still exhibit translationese features.

Some post-editing models do not rely on the source translation or human editing data (Simard et al., 2007). For instance, Freitag et al. (2019) trained a post-editor solely on monolingual data by reconstructing the original text given its round-trip translation. In our work, we incorporate stronger natural language modelling into post-editing by employing LLMs. Other translation refinement research includes combining statistical and neural systems

(Novak et al., 2016; Niehues et al., 2016), merging APE into the NMT framework (Pal et al., 2020; Chen et al., 2022), and debiasing translationese in the latent embedding space (Dutta Chowdhury et al., 2022). The iterative editing mechanism mostly lies in non-autoregressive translation, where each output token is independent of other target positions and iterative decoding enhances output quality (Lee et al., 2018; Gu et al., 2019; Xu and Carpuat, 2021).

6.2 Translation prompting with large language models

Large language models have recently become highly effective tools for various NLP tasks (Radford et al., 2019; Brown et al., 2020; Chowdhury et al., 2022; Ouyang et al., 2022). Nowadays, optimising LLMs directly for specific tasks becomes less important since they generalize to downstream tasks even without explicit supervision. With more parameters and training data, LLMs may offer stronger performance than dedicated translation or post-editing models. The method we use to elicit a response from GPT is zero-shot prompting (Brown et al., 2020), which means affixing a description to the original task input to form a query to the model. Researchers have benchmarked LLMs’ capability to translate (Vilar et al., 2023; Zhang et al., 2023; Jiao et al., 2023; Hendy et al., 2023), and to interpret translation quality (Kocmi and Federmann, 2023; Lu et al., 2023; Xu et al., 2023).

Among the recent papers on LLM translation prompting, we identify the following to be most relevant to us. Previous findings show that GPT produces less literal translations, especially for out-of-English translations (Raunak et al., 2023a), which to some extent stands in contrast with our later human evaluation results on naturalness and fluency.

Raunak et al. (2023b) formalized post-editing as a chain-of-thought process (Wei et al., 2022) with GPT-4 and achieved promising results. Different from their focus, our work features the iterative refinement process as a means to enhance naturalness and fluency. Our work reveals that iterated refinement is better than one-off editing. The observed improvement, especially for into-English, may be attributed to the abundant English pre-training data available for LLMs. To the best of our knowledge, although the concept of iterative refinement is not new, ours is the pioneering paper in applying such strategies to LLMs for translation.

7 Conclusion and Future Work

We presented a simple way to leverage an LLM for translation refinement, which greatly helps fluency and naturalness. It is shown that our method maintains translation quality and introduces lexical and structural changes, especially for high-resource into-English translation. We have also discussed the potential of using our work to obtain diverse, fluent translations that are less translationese, as well as the limitation in automatic metrics to measure this.

On a broader note, this work connects to the concept of using LLMs to imitate collaborative translation refinement. Yet, it is important to acknowledge the high cost of running a multi-round LLM refinement. Future work can explore sentence-level refinement decisions to reduce cost.

Acknowledgement

We express our gratitude to the reviewers of this paper for their detailed and invaluable feedback and suggestions. The work also benefited from discussions with Nikolay Bogoychev and Biao Zhang. We are grateful to Laurie Burchell, Ziqin Fang, Matthias Lindemann, and Jonas Waldendorf for their participation in the human evaluation.

This work is funded by UK Research and Innovation (UKRI) under the UK government’s Horizon Europe funding guarantee [grant number 10052546].

References

Agrawal, Sweta, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2023. In-context examples selection for machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023*.

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations*.

Baker, Mona, 1996. *Corpus-based Translation Studies: The Challenges that Lie Ahead*. Benjamins Translation Library. John Benjamins Publishing Company.

Bizzoni, Yuri, Tom S Jurek, Cristina España-Bonet, Koel Dutta Chowdhury, Josef van Genabith, and Elke Teich. 2020. How human is machine translationese? comparing human and machine translations of text and speech. In *Proceedings of the 17th International Conference on Spoken Language Translation*.

Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*.

Chatterjee, Rajen, Matteo Negri, Raphael Rubino, and Marco Turchi. 2018. Findings of the WMT 2018 shared task on automatic post-editing. In *Proceedings of the Third Conference on Machine Translation*.

Chen, Pinzhen, Jindřich Helcl, Ulrich Germann, Laurie Burchell, Nikolay Bogoychev, Antonio Valerio Miceli Barone, Jonas Waldendorf, Alexandra Birch, and Kenneth Heafield. 2021. The University of Edinburgh’s English-German and English-Hausa submissions to the WMT21 news translation task. In *Proceedings of the Sixth Conference on Machine Translation*.

Chen, Kehai, Masao Utiyama, Eiichiro Sumita, Rui Wang, and Min Zhang. 2022. Synchronous refinement for neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2022*.

Chollampatt, Shamil, Raymond Hendy Susanto, Liling Tan, and Ewa Szymanska. 2020. Can automatic post-editing improve NMT? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*.

Chowdhury, Aakanksha, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. PaLM: Scaling language modeling with pathways. *arXiv preprint*.

Dutta Chowdhury, Koel, Richa Jalota, Cristina España-Bonet, and Josef Genabith. 2022. Towards debiasing translation artifacts. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.

Farhad, Akhbardeh, Arkhangorodsky Arkady, Biesialska Magdalena, Bojar Ondřej, Chatterjee Rajen, Chaudhary Vishrav, Marta R Costa-jussa, España-Bonet Cristina, Fan Angela, Federmann Christian,

- et al. 2021. Findings of the 2021 conference on machine translation (WMT21). In *Proceedings of the Sixth Conference on Machine Translation*.
- Freitag, Markus, Isaac Caswell, and Scott Roy. 2019. APE at scale and its implications on MT evaluation biases. In *Proceedings of the Fourth Conference on Machine Translation*.
- Freitag, Markus, Ricardo Rei, Nitika Mathur, Chikiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation*.
- Gellerstam, Martin. 1986. Translationese in Swedish novels translated from English. In *Translation studies in Scandinavia: Proceedings from the Scandinavian Symposium on Translation Theory II*. CWK Gleerup.
- Gu, Jiatao, Changan Wang, and Junbo Zhao. 2019. Levenshtein transformer. In *Advances in Neural Information Processing Systems*.
- Hendy, Amr, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? a comprehensive evaluation. *arXiv preprint*.
- Ive, Julia, Lucia Specia, Sara Szoc, Tom Vanallemeersch, Joachim Van den Bogaert, Eduardo Farah, Christine Maroti, Artur Ventura, and Maxim Khalilov. 2020. A post-editing dataset in the legal domain: Do we underestimate neural machine translation quality? In *Proceedings of the Twelfth Language Resources and Evaluation Conference*.
- Jiao, Wenxiang, Wenxuan Wang, Jen tse Huang, Xing Wang, and Zhaopeng Tu. 2023. Is ChatGPT a good translator? Yes with GPT-4 as the engine. *arXiv preprint*.
- Junczys-Dowmunt, Marcin and Roman Grundkiewicz. 2018. MS-UEdin submission to the WMT2018 APE shared task: Dual-source transformer for automatic post-editing. In *Proceedings of the Third Conference on Machine Translation*.
- Kaplan, Jared, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint*.
- Knight, Kevin and Ishwar Chander. 1994. Automated postediting of documents. In *Proceedings of the Twelfth AAAI National Conference on Artificial Intelligence*.
- Kocmi, Tom and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. *arXiv preprint*.
- Kocmi, Tom, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, et al. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation*.
- Koppel, Moshe and Noam Ordan. 2011. Translationese and its dialects. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*.
- Lee, Jason, Elman Mansimov, and Kyunghyun Cho. 2018. Deterministic non-autoregressive neural sequence modeling by iterative refinement. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
- Lembersky, Gennadi, Noam Ordan, and Shuly Wintner. 2012. Language Models for Machine Translation: Original vs. Translated Texts. *Computational Linguistics*.
- Lu, Qingyu, Baopu Qiu, Liang Ding, Liping Xie, and Dacheng Tao. 2023. Error analysis prompting enables human-like translation evaluation in large language models: A case study on ChatGPT. *arXiv preprint*.
- Niehues, Jan, Eunah Cho, Thanh-Le Ha, and Alex Waibel. 2016. Pre-translation for neural machine translation. In *Proceedings of the 26th International Conference on Computational Linguistics*.
- Novak, Roman, Michael Auli, and David Grangier. 2016. Iterative refinement for machine translation. *arXiv preprint*.
- Ouyang, Long, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*.
- Pal, Santanu, Hongfei Xu, Nico Herbig, Sudip Kumar Naskar, Antonio Krüger, and Josef van Genabith. 2020. The transference architecture for automatic post-editing. In *Proceedings of the 28th International Conference on Computational Linguistics*.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*.
- Popović, Maja. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*.
- Radford, Alec, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *openai.com*.

- Raunak, Vikas, Arul Menezes, Matt Post, and Hany Hassan. 2023a. Do GPTs produce less literal translations? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*.
- Raunak, Vikas, Amr Sharaf, Hany Hassan Awadallah, and Arul Menezes. 2023b. Leveraging GPT-4 for automatic translation post-editing. *arXiv preprint*.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*.
- Simard, Michel, Cyril Goutte, and Pierre Isabelle. 2007. Statistical phrase-based post-editing. In *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics*.
- Teich, Elke. 2003. *Cross-Linguistic Variation in System and Text: A Methodology for the Investigation of Translations and Comparable Texts*. De Gruyter Mouton.
- Toral, Antonio. 2019. Post-editese: An exacerbated translationese. In *Proceedings of Machine Translation Summit XVII*.
- Tran, Chau, Shruti Bhosale, James Cross, Philipp Koehn, Sergey Edunov, and Angela Fan. 2021. Facebook AI’s WMT21 news translation task submission. In *Proceedings of the Sixth Conference on Machine Translation*.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*.
- Vilar, David, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting PaLM for translation: Assessing strategies and performance. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*.
- Wang, Longyue, Mu Li, Fangxu Liu, Shuming Shi, Zhaopeng Tu, Xing Wang, Shuangzhi Wu, Jiali Zeng, and Wen Zhang. 2021. Tencent translation system for the WMT21 news translation task. In *Proceedings of the Sixth Conference on Machine Translation*.
- Wei, Daimeng, Zongyao Li, Zhanglin Wu, Zhengzhe Yu, Xiaoyu Chen, Hengchao Shang, Jiaxin Guo, Minghan Wang, Lizhi Lei, Min Zhang, et al. 2021. HWTSC’s participation in the WMT 2021 news translation shared task. In *Proceedings of the Sixth Conference on Machine Translation*.
- Wei, Jason, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*.
- Xu, Weijia and Marine Carpuat. 2021. EDITOR: An edit-based transformer with repositioning for neural machine translation with soft lexical constraints. *Transactions of the Association for Computational Linguistics*.
- Xu, Wenda, Danqing Wang, Liangming Pan, Zhenqiao Song, Markus Freitag, William Yang Wang, and Lei Li. 2023. INSTRUCTSCORE: Towards explainable text generation evaluation with automatic feedback. *arXiv preprint*.
- Zhang, Biao, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study. In *Proceedings of the 40th International Conference on Machine Learning*.

Detector–Corrector: Edit-Based Automatic Post Editing for Human Post Editing

Hiroyuki Deguchi¹ Masaaki Nagata² Taro Watanabe¹

¹Nara Institute of Science and Technology

²NTT Communication Science Laboratories, NTT Corporation

¹{deguchi.hiroyuki.db0, taro}@is.naist.jp

²masaaki.nagata@ntt.com

Abstract

Post-editing is crucial in the real world because neural machine translation (NMT) sometimes makes errors. Automatic post-editing (APE) attempts to correct the outputs of an MT model for better translation quality. However, many APE models are based on sequence generation, and thus their decisions are harder to interpret for actual users. In this paper, we propose “detector–corrector”, an edit-based post-editing model, which breaks the editing process into two steps, error detection and error correction. The detector model tags each MT output token whether it should be corrected and/or reordered while the corrector model generates corrected words for the spans identified as errors by the detector. Experiments on the WMT’20 English–German and English–Chinese APE tasks showed that our detector–corrector improved the translation edit rate (TER) compared to the previous edit-based model and a black-box sequence-to-sequence APE model, in addition, our model is more explainable because it is based on edit operations.

1 Introduction

Neural machine translation (NMT) (Sutskever et al., 2014; Bahdanau et al., 2015; Luong et al., 2015; Wu et al., 2016; Vaswani et al., 2017) sometimes make errors (Ott et al., 2018), and post-editing is crucial in the real world to correct the

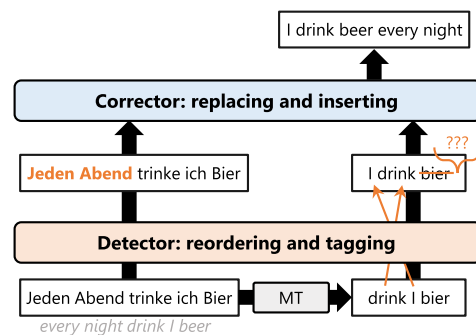


Figure 1: Overview of the post-editing process of our detector–corrector model. The detector tags as “Jeden Abend” is untranslated, “drink” and “I” should be reordered, etc. The corrector generates the word sequence for replacement and insertion.

mis-translations. Automatic post-editing (APE) attempts to correct and refine the translations generated by MT models (MT sentences) for better translation quality. However, many APE models are based on sequence generation (Junczys-Dowmunt and Grundkiewicz, 2018; Correia and Martins, 2019; Sharma et al., 2021; Chatterjee et al., 2019; Chatterjee et al., 2020; Bhattacharyya et al., 2022), and their decision for correction is harder to interpret due to the black-box nature of the generation models.

Some prior work (Malmi et al., 2019; Gu et al., 2019; Omelianchuk et al., 2020; Stahlberg and Kumar, 2020; Mallinson et al., 2020; Mallinson et al., 2022) showed that edit-based models improve interpretability in monolingual text editing, e.g., grammatical error correction (GEC), compared with sequence-to-sequence models. The APE task can be regarded as a text edit task in terms of rewriting MT sentences, but differs from general monolingual text editing tasks in that it uses cross-lingual information from source sentences, such as inserting untranslated words and re-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

ordering translation words. For example, if an edit-based model cannot perform reordering, it is represented as deletion and insertion, which increases the number of edit operations and makes it harder for humans to interpret the edit.

In this paper, we propose “detector–corrector”, an edit-based post-editing model, in which the post-editing process is broken into two steps for assisting human post-editing: error detection and error correction. We designed our model after interviewing with professional translators regarding the post-editing process; specifically, they first spot errors and then make corrections, and omission errors are crucial for the editing process. The overview of our detector–corrector model is shown in Figure 1. The detector model, which extends a word-level quality estimation (QE) model, tags each MT output token as whether it should be corrected and/or reordered and identifies which source tokens are not translated in the MT sentence. Then, the corrector model receives the annotated source and MT sentences and corrects words for each span identified as incorrect in the detector model. Our corrector model can insert any number of spans of variable length. In addition, we propose data augmentation methods especially designed for the detector and corrector models to enhance each model, and lightweight iterative refinement to improve the inference speed.

Experiments on the WMT’20 English–German (En–De) and English–Chinese (En–Zh) APE tasks showed that our detector–corrector improved translation edit rate (TER) (Snover et al., 2006) compared to not only an edit-based model (Gu et al., 2019) but also a black-box sequence-to-sequence model by 0.7 points in En–De and En–Zh. Moreover, our model is more explainable than sequence-to-sequence models because it is based on edit operations and it can be integrated into computer-aided translation tools (Herbig et al., 2020).

2 Background and Related Work

2.1 Edit-Based Model

Chen et al. (2020) have built an edit-based GEC system that detects erroneous spans and then corrects the words within the detected erroneous spans. GECToR (Omelianchuk et al., 2020) is also an edit-based GEC mode, in which the model predicts the error type tag for each word, and then words identified as errors are corrected according

to the rules for each tag type.

Levenshtein Transformer (Gu et al., 2019), a non-autoregressive Transformer encoder-decoder model, predicts deletion, placeholder insertion, and word filling. It can be used for the APE task by rewriting an MT sentence, but it cannot represent reordering and detecting untranslated words. Seq2Edits (Stahlberg and Kumar, 2020) edits an input text by span tagging and replacement prediction to improve interpretability for text-editing tasks. However, it is not suitable for the APE task because it only monotonically edits an MT output from left to right according to the tags and cannot perform reordering of spans or inserting missing words which often occur in erroneous translations. FELIX (Mallinson et al., 2020) breaks down text editing into three components: tagging, reordering, and word in-filling. It performs tagging using a pre-trained encoder model like BERT, reordering using a pointer network, and predicting words of replacement and insertion using a masked language model. However, it does not explicitly use source information. In addition, word insertion is predicted non-autoregressively; thus, the number of words to be inserted must be given in advance for the insertion operation, which is not trivial. Edit5 (Mallinson et al., 2022) uses the T5 (Rafael et al., 2020) encoder-decoder and decomposes the editing process into (1) tagging that decides which tokens are kept, (2) reordering the input tokens, and (3) insertion that infills the missing tokens. Unlike FELIX, Edit5 uses the autoregressive T5 decoder for word prediction, allowing for variable length insertion. However, the positions that can be inserted depend on the special tokens used in pre-training of T5 for filling masked spans, e.g., `<extra_id.6>` as `<pos6>`; thus, the number of positions that can be inserted is limited to those observed in pre-training.

2.2 Word-Level Quality Estimation

The word-level quality estimation task estimates the word-level quality of MT sentences, which is closely related to the post-editing task. It is divided into three binary classifications (Specia et al., 2020): MT-tag, MT-gap, and SRC-tag. MT-tag detects erroneous words in MT sentences. MT-gap predicts where to insert untranslated words in MT sentences, and SRC-tag detects untranslated source words.

Predictor-estimator model (Kim et al., 2017a;

Kim et al., 2017b) is a well-known architecture for the word-level quality estimation task, in which the predictor is used for feature extraction from translation results while the estimator estimates the translation quality based on the features from the predictor. Ding et al. (2021) used Levenshtein Transformer (Gu et al., 2019) for the word-level quality estimation task. Their method uses the edit probabilities of deletion and insertion of Levenshtein Transformer as tag prediction probabilities instead of explicitly predicting OK/BAD tags. DirectQE (Cui et al., 2021) is a pre-training method designed for the QE task, which consists of two components: generator and detector. In pre-training, The generator rewrites words by a cross-lingual masked language model, then the detector detects the replaced words. After pre-training, the detector model is fine-tuned with real QE data. SiameseTransQuest (Ranasinghe et al., 2020) employed the word-level QE architecture using XLM-R for the sentence-level quality estimation task, and they showed that using XLM-R is effective in the QE task. Ranasinghe et al. (2021) demonstrated that the fine-tuned XLM-R predicts word-level QE on other language pairs than a language pair that is trained explicitly, i.e., the model can perform zero-shot QE.

2.3 Automatic Post Editing

The automatic post-editing (APE) task aims to improve the translation quality by editing translations generated from black-box MT models (Chatterjee et al., 2020). The APE system receives the source and MT sentences and generates the post-edited (PE) sentence. This task mainly evaluates correction performance using translation edit rate (TER) (Snover et al., 2006) based on the edit distance between the human-revised translation and the corrected sentence.

Correia and Martins (2019) built a sequence-to-sequence APE system by only fine-tuning pre-trained BERT models, in which weight initialization is carefully designed to employ pre-trained weights for both encoder and decoder. In the APE shared task, the high-ranked systems often employ Transformer encoder-decoder architectures with pre-trained models (Chatterjee et al., 2020; Bhattacharyya et al., 2022; Yang et al., 2020; Wang et al., 2020; Lee et al., 2020; Deoghare and Bhattacharyya, 2022; Huang et al., 2022). The sequence-to-sequence model, which

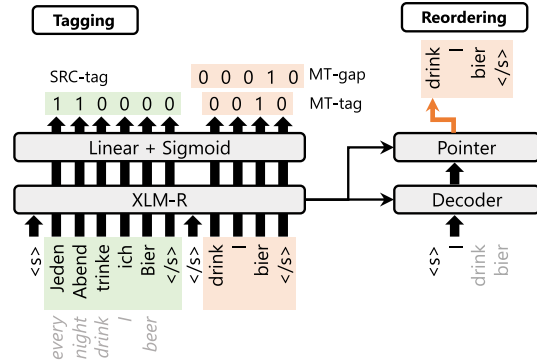


Figure 2: Overview of our detector model. The model detects OK and BAD tags as 0 and 1, respectively.

learns post-editing in an end-to-end manner, can achieve high translation quality; however, it cannot explicitly expose the editing process, making it hard to utilize the model in scenarios that require manual checking. The copy mechanism (Gu et al., 2016) can be used for APE tasks by copying words in MT sentences that do not need to be modified (Huang et al., 2019). This model can show us edited and non-edited words using the copy probability. Neural Programmer-Interpreter (NPI) (Vu and Haffari, 2018) generates PE sentences by predicting the edit actions and the target tokens comprising three editing operations: keep, delete, and insert. Although NPI is more interpretable than the sequence-to-sequence models, it cannot represent reordering nor differentiate replacement and insertion. Deoghare et al. (2023) incorporated the word-level quality estimation into an APE model. Their model predicts which word should be edited through multi-task learning; however, it cannot use human-annotated QE tags because the information of QE tags, which is passed to the decoder, is represented as hidden vectors.

3 Proposed Model: Detector-Corrector

3.1 Detector

Our detector model (Figure 2) predicts shift and edit operations based on translation edit rate (TER) (Snover et al., 2006). TER iteratively reorders an input sequence to minimize the edit distance from the target sequence, called “shift” operation, then calculates edit distance between the reordered input sequence and the target sequence, called “edit” operations. To represent this TER behavior, our detector model performs tagging to predict whether edits are needed (“Tagging” in Figure 2), and reordering of the given

MT sentence with a pointer network (Vinyals et al., 2015) (“Reordering” in Figure 2). Let $\mathbf{x} = (x_1, \dots, x_{|\mathbf{x}|}) \in \mathcal{V}^*$ and $\mathbf{y} = (y_1, \dots, y_{|\mathbf{y}|}) \in \mathcal{V}^*$ denote the given source sentence and its translation generated by machine translation (MT sentence), respectively, where \mathcal{V}^* is the Kleene closure of the vocabulary¹ \mathcal{V} . Note that both \mathbf{x} and \mathbf{y} always have the end-of-sentence symbol “</s>” as the last tokens, i.e., $x_{|\mathbf{x}|} = y_{|\mathbf{y}|} = \text{“</s>”}$. Let $\mathbf{x} \circ \mathbf{y}$ be the concatenated sequence, where \circ represents the join operation with a separator token between the sequences². XLM-RoBERTa (XLM-R) encoder (Conneau et al., 2020) encodes the concatenated sequence $\mathbf{x} \circ \mathbf{y}$ into D -dimensional hidden vectors through L layers $\mathbf{H}^{(L)} = (\mathbf{h}_1^{(L)}, \dots, \mathbf{h}_{|\mathbf{x} \circ \mathbf{y}|}^{(L)})^\top \in \mathbb{R}^{|\mathbf{x} \circ \mathbf{y}| \times D}$.

Tagging To perform tagging, we train a word-level quality estimation model. In particular, the detector model performs three binary classifications as defined by Specia et al. (2020): MT-tag, MT-gap, and SRC-tag.

Let $\mathbf{o}^T \in \{0, 1\}^{|\mathbf{y}|}$ denote the MT-tag which represents whether an MT token would be edited, i.e., $o_i^T = 1$ if y_i is deletion or replacement in a TER edit sequence, e.g., “bier” in Figure 2. The MT-tag classification identifies whether an MT token should be edited based on the bad probabilities:

$$p_i^T := p(o_i^T = 1 | \mathbf{x}, \mathbf{y}) = \sigma(\mathbf{w}_T^\top \mathbf{h}_{y_i}^{(l_T)}), \quad (1)$$

where $\mathbf{w}_T \in \mathbb{R}^D$ is a learned parameter for MT-tag prediction, $1 \leq l_T \leq L$ denotes the layer used for MT-tag prediction, and $\sigma : \mathbb{R} \rightarrow [0, 1]$ is a sigmoid function. Note that $\mathbf{h}_{y_i}^{(l)}$ is a row of $\mathbf{H}^{(l)}$, which is the hidden vector corresponding to the token y_i in the l -th layer.

Similarly, MT-gap classification predicts whether some words need to be inserted at a token boundary in the MT sentence based on the insertion probabilities:

$$p_i^G := p(o_i^G = 1 | \mathbf{x}, \mathbf{y}) = \sigma(\mathbf{w}_G^\top [\mathbf{h}_{y_{i-1}}^{(l_G)}; \mathbf{h}_{y_i}^{(l_G)}]), \quad (2)$$

where $\mathbf{o}^G \in \{0, 1\}^{|\mathbf{y}|}$ represents insertion in a TER edit sequence, e.g., the token boundary between

¹We employ XLM-R, a multilingual encoder; thus, the vocabulary is shared between the source and target languages.

²In XLM-R, the class token is represented by “<s>”, and two sentences are joined by “</s>” symbols, like “<s> a b c </s> </s> A B </s>”. We regard the first symbol as the end-of-sentence symbol of the first sentence, i.e., $x_{|\mathbf{x}|}$, and the second one as the separator token.

“bier” and “</s>” in Figure 2. $\mathbf{w}_G \in \mathbb{R}^{2D}$ is a learned parameter for MT-gap prediction, $1 \leq l_G \leq L$ denotes the layer used for MT-gap prediction, and $[\cdot; \cdot]$ denotes the concatenation of two vectors. Note that y_0 is the separator token between the source and MT sentences.

Likewise, the SRC-tag $\mathbf{o}^S \in \{0, 1\}^{|\mathbf{x}|}$ is constructed from a source-target word alignment as $x_i = 1$ if x_i is not aligned to any target token like “Jeden” and “Abend” in Figure 2. In this paper, we used AWESOME-ALIGN (Dou and Neubig, 2021) to obtain the gold alignment. The SRC-tag classification predicts whether a source token is untranslated or not using the probabilities:

$$p_i^S := p(o_i^S = 1 | \mathbf{x}, \mathbf{y}) = \sigma(\mathbf{w}_S^\top \mathbf{h}_{x_i}^{(l_S)}), \quad (3)$$

where $\mathbf{w}_S \in \mathbb{R}^D$ is a learned parameter for SRC-tag prediction and $1 \leq l_S \leq L$ denotes the layer used for SRC-tag prediction.

During inference, each tag \mathbf{o}^T , \mathbf{o}^G , and \mathbf{o}^S are respectively predicted to be “BAD” when each probability p_i is greater than 0.5, and “OK” otherwise.

Reordering Our detector also predicts reordering by generating the reordered sequence $\bar{\mathbf{y}} = (\bar{y}_1, \dots, \bar{y}_{|\bar{\mathbf{y}}|})$ using the pointer network (Vinyals et al., 2015) at the top of the decoder. It autoregressively selects the next token for each timestep from the MT sentence according to the probability p^R , as follows:

$$\bar{\mathbf{y}}^* = \operatorname{argmax}_{(\bar{y}_1, \dots, \bar{y}_{|\bar{\mathbf{y}}|})} \prod_{i=1}^{|\bar{\mathbf{y}}|} p^R(\bar{y}_i | \mathbf{x}, \mathbf{y}, \bar{\mathbf{y}}_{<i}), \quad (4)$$

$$p^R(\bar{y}_i = y_j | \mathbf{x}, \mathbf{y}, \bar{\mathbf{y}}_{<i}) \propto \exp(\mathbf{k}_{y_j}^\top \mathbf{q}_{\bar{y}_i}), \quad (5)$$

$$\mathbf{k}_{y_j} = \mathbf{W}_k \mathbf{h}_{y_j}, \quad (6)$$

$$\mathbf{q}_{\bar{y}_i} = \mathbf{W}_q \text{Decoder}(\bar{\mathbf{y}}_{<i}, \mathbf{H}^{(L)}), \quad (7)$$

where $\text{Decoder} : \mathcal{V}^* \times \mathbb{R}^{|\mathbf{x} \circ \mathbf{y}| \times D} \rightarrow \mathbb{R}^D$ is a Transformer decoder that computes a hidden vector of the i -th step $\mathbf{q}_{\bar{y}_i}$ from the given encoder hidden vectors and the prefix of reordered sequence. $\mathbf{W}_q \in \mathbb{R}^{D \times D}$ and $\mathbf{W}_k \in \mathbb{R}^{D \times D}$ are the learned parameters, and $\bar{\mathbf{y}}^*$ is the reordered sequence predicted by the model. Note that the hidden vectors $\mathbf{H}^{(L)}$ are computed using the same encoder as used in tagging.

During inference, the tokens of the MT sentence and their corresponding MT-tag and MT-gap are reordered according to the order of $\bar{\mathbf{y}}^*$. Note that

the MT-gap tags are reordered in accordance with the order of their right-side tokens of boundaries. For example, in Figure 2, the MT-gap model predicts that some words need to be inserted at the token boundary between “bier” and “</s>”, and the boundary position is attached to the left of “</s>” after reordering.

Objective function We trained the MT-tag, MT-gap, and SRC-tag classifications by minimizing their objective functions, \mathcal{L}_T , \mathcal{L}_G , and \mathcal{L}_S , computed by the binary cross-entropy, as follows:

$$-\sum_i (o_i \log p_i + (1 - o_i) \log(1 - p_i)), \quad (8)$$

where $o_i \in \{0, 1\}$ is the ground truth label of the probability p_i . The model is also trained to generate reordered MT sentences by minimizing the following cross-entropy:

$$\mathcal{L}_R = -\sum_{i=1}^{|\mathbf{y}|} \log p^R(\bar{y}_i | \mathbf{x}, \mathbf{y}, \bar{\mathbf{y}}_{<i}), \quad (9)$$

where the gold reordered sequence is created from the TER shift alignment. Finally, our detector model is trained by minimizing the following objective \mathcal{L} through multi-task learning:

$$\mathcal{L} = \mathcal{L}_T + \mathcal{L}_G + \mathcal{L}_S + \mathcal{L}_R. \quad (10)$$

Note that all loss functions in \mathcal{L} are computed during a single forward pass since the encoder parameters are shared between all tagging and reordering predictions.

3.2 Corrector

The corrector model (Figure 3) corrects the reordered MT sentence by generating tokens corresponding to the erroneous spans identified by MT-tag and MT-gap predictions. The corrector represents edit operations by predicting zero words in a bad span for deletion, one or more words in a bad span for replacement, and one or more words in an insertion span for insertion, as shown on the output of the decoder in Figure 3.

First, the tags predicted by the detector model are used to annotate the source sentence and its corresponding reordered MT output as span tags. In the source sentence, <bad> and </bad> tags are inserted to the beginning and end of untranslated spans, respectively, using the SRC-tag \mathbf{o}^S , as shown on the left side of the input of the XLM-R encoder in Figure 3. Similarly, <bad> and

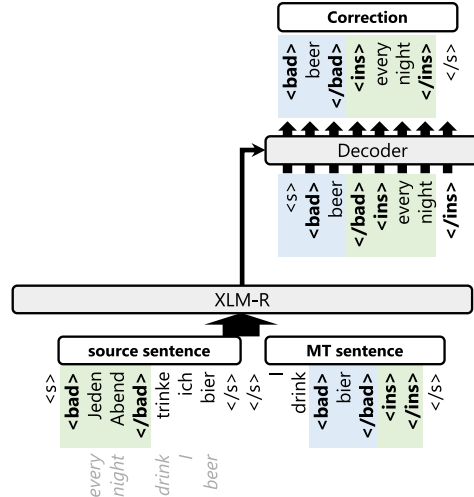


Figure 3: Token generation within each tagged span by our corrector model.

</bad> tags are inserted into reordered MT output where identified by the MT-tag tagging \mathbf{o}^T in addition to the <ins> and </ins> tags to the positions that need to be inserted words, as shown on the right side of the input of the XLM-R encoder in Figure 3.

Next, the annotated source and reordered MT sentences are concatenated with the separator token and fed into the encoder. We initialize the corrector encoder with XLM-R as well as the detector model in order to preserve consistency with the subword unit tags used in the detector. Then, the decoder generates tokens for all tagged spans in the left-to-right manner until the number of corrected spans satisfies the number of bad and insertion spans in the annotated reordered MT sentence. Finally, our detector–corrector outputs a corrected target sentence by replacing each tagged span of the MT sentence with a token sequence predicted by the corrector decoder.

Our corrector can be regarded as a translation suggestion (TS) model (Yang et al., 2022a; Yang et al., 2022b), in which better alternative translations are suggested phrase-by-phrase by replacing incorrect translation spans. Our model differs from TS models in that untranslated spans in source sentences are explicitly identified and incorrect translations and/or insertions are clearly differentiated by the bad and insertion tags, respectively. Furthermore, MT sentences are reordered and multiple spans are corrected in our model, which are out of the scope of the TS task³.

³The TS task assumes only a single incorrect span for each sentence and does not treat reordering.

3.3 Data Augmentation

3.3.1 Data Augmentation for Detector

Since the detector–corrector is trained to correct only erroneous spans identified by the detector, improving the tagging accuracy will directly lead to improved translation quality. For this purpose, we create the synthetic data from the reference translations of the training data and let the detector learn the editing operations of deletion, replacement, and insertion. We randomly delete tokens with a probability of 5%, insert tokens with a probability of 10%, and replace tokens with a probability of 30%. We employ XLM-R to fill the masked tokens for the replacement and insertion decision.

3.3.2 Data Augmentation for Corrector

The training data for the corrector model is created from the tokens for each span identified as an error using the oracle annotated source and MT sentences. However, the detector might make wrong decision during inference, which might cause a large discrepancy between the training and inference for the corrector. In addition, the performance of the corrector might suffer from the limited coverage of the vocabulary in the training data when compared with a conventional sequence-to-sequence MT model. For these reasons, we employ two simple data augmentation methods for the corrector model without additional computational cost: MT training and PE training. These two augmentation methods are orthogonal with each other; thus, they can be combined.

MT Training In MT training, the corrector model is trained to predict the PE sentence from only the source sentence without the corresponding MT sentence. To preserve the model consistency, an MT output is treated as an empty text by augmenting with “<ins> </ins>” so that the model learns to insert the whole PE sentence from the empty MT sentence. The encoder input sequence of MT training is formulated as follows:

$$\langle \text{bad} \rangle \mathbf{x} \langle \text{/bad} \rangle \circ \langle \text{ins} \rangle \langle \text{/ins} \rangle, \quad (11)$$

and the corrector is trained to generate the post-edited sentence with the insertion, i.e., $\langle \text{ins} \rangle \mathbf{y}^{\text{PE}} \langle \text{/ins} \rangle$, where $\mathbf{y}^{\text{PE}} \in \mathcal{V}^*$ is the post-edited sentence.

PE Training PE training differs from MT training in that the MT sentences are given. The corrector model is trained to generate the whole PE

sentence from the given source and MT sentences. This is the same setting as the standard sequence-to-sequence APE model training, except that the MT sentence is explicitly annotated as “<bad>”. To maintain model consistency, the whole MT sentence is treated as a bad span to be corrected:

$$\mathbf{x} \circ \langle \text{bad} \rangle \mathbf{y} \langle \text{/bad} \rangle, \quad (12)$$

and the model learns to replace the MT sentence with the PE sentence, i.e., the model is trained to generate $\langle \text{bad} \rangle \mathbf{y}^{\text{PE}} \langle \text{/bad} \rangle$.

3.4 Lightweight Iterative Refinement

The detector model detects each erroneous span in a non-autoregressive manner; thus, a single inference may not generate sufficiently correct PE sentences that are consistent across the entire sentence. To address such issues, some prior non-autoregressive models (Gu et al., 2019; Kasai et al., 2020; Omelianchuk et al., 2020) decode sequences by iteratively feeding the output into the model. We follow the practice by iteratively refining an MT sentence by treating the post-edited sentence corrected by our model as an MT output, i.e., the corrected sentence in the $k - 1$ -th iteration is used as the input of the detector model in the k -th iteration. However, the iterative refinement approach demands huge computation in particular for our approach, in which an end-to-end inference predicts three edit operations in the following order: tagging, reordering, and correcting.

Tagging can be predicted with only a single forward pass of the detector encoder, and correcting can be finished very quickly since it generates only a few words for each erroneous span. In contrast, reordering is relatively slower than the other operations because the decoder runs for the length of the MT sentence in an auto-regressive manner.

In order to overcome such bottleneck, we propose lightweight refinement, in which inference is carried out only by predicting tags and generating correct tokens without reordering after the second time in the iterative refinement.

4 Experiments

4.1 Setup

We compared the translation quality of our detector–corrector with that of the sequence-to-sequence (seq2seq) APE model and Levenshtein Transformer (LevT) (Gu et al., 2019). We evaluated TER (\downarrow T), BLEU (\uparrow B), and COMET (\uparrow C) us-

ing SACREBLEU (Post, 2018) and COMET⁴ (Rei et al., 2020; Rei et al., 2022) in the WMT’20 English–German (En–De) and English–Chinese (En–Zh) automatic post-editing tasks.

Datasets Training data came from WMT’20 APE tasks, which were created from wikipedia articles that contain 7,000 sentences, and we applied upsampling by 20 times to them. In addition to the provided data, we created additional training data that consists of ⟨source sentence, MT sentence, PE sentence⟩ triplets using a parallel corpus following the idea from Negri et al. (2018). In particular, we randomly sampled 2 million sentences from the training data of the WMT’19 En–De and En–Zh translation tasks and translated them with MT models, which were used to generate the data for the APE tasks (Fomicheva et al., 2020). As described in Section 3.3, the training data for the detector and corrector were further augmented. The data statistics are shown in the appendix (Table 10).

Models The seq2seq APE model, LevT, and our detector–corrector comprise the XLM-R large encoder and Transformer decoder. The seq2seq, LevT, and corrector models were trained in 60,000 steps, and the detector model was trained in 40,000 steps. All models were optimized by Adam optimizer ($\beta_1 = 0.9, \beta_2 = 0.98$). The learning rate was linearly increased up to 4,000 steps and then decayed proportional to the inverse square root of the training steps. The beam size was set to 5, and the length penalty was set to $\alpha = 1.0$. We saved checkpoints of all models for every 1,000 steps and took an average of the last 5 checkpoints. The LevT edited the MT sentences 5 times iteratively, and the detector–corrector edited 4 times, i.e., $k = 4$, by tuning on the development set. For tagging, we used the intermediate representations of the 20th layer, i.e., $l_T = l_G = l_S = 20$ in En–De, and the 24th layer, i.e., $l_T = l_G = l_S = 24$ in En–Zh. The details of each model are shown in the appendix (Table 9).

4.2 Results

Our main results are shown in Table 1. Our detector–corrector model improved TER and BLEU from both LevT and seq2seq models. Especially in TER, detector–corrector outperforms the

Dataset	Model	↓T	↑B	↑C
En–De	do nothing (MT)	31.3	50.2	77.1
	seq2seq	28.4	53.3	77.7
	LevT (Gu et al., 2019)	31.9	49.4	75.6
	detector–corrector	27.7 [†]	53.6	79.6 [†]
En–Zh	do nothing (MT)	58.3	24.3	86.3
	seq2seq	56.7	26.0	89.4 [†]
	LevT (Gu et al., 2019)	59.3	23.6	86.0
	detector–corrector	56.0	26.1	89.2

Table 1: Comparison of post-editing performance in the WMT’20 En–De and En–Zh APE tasks. Do nothing (MT) does not edit MT sentences and the scores are calculated between MT and PE sentences. The best scores of each dataset are emphasized by the **bold** font. The symbol † indicates that the score difference is statistically significant ($p < 0.05$) between seq2seq and detector–corrector.

Model	En–De			En–Zh		
	↓T	↑B	↑C	↓T	↑B	↑C
ours	27.7 [†]	53.6 [†]	79.6 [†]	56.0 [†]	26.1 [†]	89.2 [†]
- light-iter	28.9	52.1	77.7	56.6	25.5	88.0
-- MT training	29.3	51.5	77.7	56.6	25.4	88.3
-- PE training	29.2	51.8	77.7	56.6	25.2	88.3
-- DAug for corrector	30.2	50.1	77.6	57.0	24.9	88.6
--- DAug for detector	31.2	49.0	77.1	61.2	22.7	86.7

Table 2: Ablation study of our methods in the WMT’20 En–De and En–Zh APE tasks. The symbol † indicates that the score difference is statistically significant ($p < 0.05$) between “ours” and “- light-iter”.

black-box seq2seq model by 0.7 % in En–De and En–Zh while providing the editing process.

Table 2 shows the ablation study of our proposed methods. In the table, “light-iter” denotes the lightweight iterative refinement, and “DAug” denotes data augmentation. The results show that both lightweight iterative refinement and data augmentation for the detector and corrector are effective, which improve the TER scores by 3.5 % in En–De and 5.2 % in En–Zh compared to the vanilla detector–corrector.

Our data augmentation for the detector can be used for other baseline models, seq2seq and LevT⁵. To confirm that the data augmentation is effective for our model, we also trained the baseline models using the augmented data. Table 3 shows that the translation quality of baseline models trained on the augmented data. Unlike the “DAug for detector” row in Table 2, there is no improvement in all metrics of more than 1 % even if the augmented data is used. This is because the

⁴<https://huggingface.co/Unbabel/wmt22-comet-da>

⁵The data augmentation for corrector cannot be applied to other models because they have been already trained to generate the whole target sentence.

Dataset	Model	↓T		↑B		↑C	
		w/o	w	w/o	w	w/o	w
En-De	seq2seq	28.4	28.4	53.3	52.9	77.7	78.0
	LevT	31.9	32.1	49.4	49.0	75.6	75.8
En-Zh	seq2seq	56.7	57.0	26.0	26.0	89.4	89.5
	LevT	59.3	59.9	23.6	23.4	86.0	86.1

Table 3: Translation quality of baseline models trained using our data augmentation for the detector.

Tagging	Dataset	DAug	MCC	F1-OK	F1-BAD
Target	En-De	w/o	0.468	0.935	0.523
		w/	0.475	0.937	0.526
	En-Zh	w/o	0.505	0.893	0.602
		w/	0.537	0.902	0.619
Source	En-De	w/o	0.782	0.985	0.794
		w/	0.791	0.985	0.805
	En-Zh	w/o	0.641	0.943	0.695
		w/	0.676	0.948	0.724

Table 4: Word-level quality estimation performance of our detector model.

data augmentation for the detector is designed to enhance word-level quality estimation.

To summarize, we confirmed that our model outperformed LevT and a black-box seq2seq model, and our approaches mitigate the translation quality degradation issue caused by predicting tags in a non-autoregressive manner and being trained from only a vocabulary limited to correction words.

5 Discussion

5.1 Accuracy of the Detector

We evaluated the tagging performance of our detector model and investigated the effectiveness of data augmentation for the detector. Since tags are predicted on subword units, we assigned a BAD tag to a word if one of the subwords in the word was assigned a BAD tag. The gold tags are calculated from the TER edit sequence after applying the shift operations in the same way as described in Section 3.1.

Table 4 shows the results of the word-level quality estimation. In the table, ‘‘MCC’’ denotes Matthews correlation coefficient (Matthews, 1975). ‘‘Target’’ and ‘‘Source’’ are the target-side tagging, i.e., MT-tag and MT-gap without distinction, and the source-side tagging, i.e., SRC-tag, respectively. We only compared our models with and without data augmentation. This is because in the

Dataset	Model	↓T	↑B	↑C
En-De	do nothing (MT)	31.3	50.2	77.1
	detector-corrector	27.7	53.6	79.6
	w/ oracle tags	13.8	74.6	82.9
		(-13.9)	(+21.0)	(+3.3)
En-Zh	do nothing (MT)	58.3	24.3	86.3
	detector-corrector	56.0	26.1	89.2
	w/ oracle tags	33.2	46.6	90.1
		(-22.8)	(+20.5)	(+0.9)

Table 5: Correction performance in the WMT’20 En-De and En-Zh APE tasks when the erroneous spans are given manually.

WMT’20 word-level QE task, the target-side tags are produced from TER edit operations without shift operations, and the source-side tags are produced by FAST_ALIGN⁶ (Dyer et al., 2013), while in our model the target-side tags include the shift operation and the source-side tags are produced by AWESOME-ALIGN. The results show that the data augmentation for the detector improved the all MCC scores, which has the direct impact to the improvements measured by BLEU and TER for our detector-corrector as shown in Table 2.

5.2 Correction Performance of Oracle Tagged Sentences

We evaluated the performance of the corrector model for oracle tags, assuming a setting in which error spans are given manually. Oracle tags were given from the TER alignment between the MT sentence and the reference translation as well as the supervision in the training data.

In Table 5, ‘‘w/ oracle tags’’ shows the result of oracle correction in the WMT’20 En-De and En-Zh APE tasks. The results showed that when given the ideal tags, the correction performance significantly improved by -13.9 and -22.8 % TER, +21.0 and +20.5 % BLEU, and +3.3 and +0.9 % COMET in En-De and En-Zh, respectively. This means that the corrector model has been successfully trained, and a further improvement in post-editing performance can be achieved by improving the accuracy of the detector model.

5.3 Ablation Study of Reordering

We also investigated the effectiveness of using the reordering operation. The training data for the model without reordering was created from the edit alignments based on the edit distance. We

⁶SIMALIGN (Jalili Sabet et al., 2020) is employed since the WMT’21 word-level QE task.

Reordering	En-De			En-Zh		
	↓T	↑B	↑C	↓T	↑B	↑C
w/	28.9	52.1	77.7	56.6	25.5	88.0
w/o	28.9	52.4	78.2	57.4	24.9	88.1

Table 6: Translation quality of detector-corrector with and without reordering. Note that we evaluated translation quality on the results of the first iteration in iterative refinement.

Reordering	En-De		En-Zh	
	# of edits	TER _{MT}	# of edits	TER _{MT}
w/	2,506	17.6	5,603	31.6
w/o	2,614	18.5	7,410	38.0

Table 7: The total number of spans tagged by the detector and TER scores that measured the amount of editing from the MT sentence to the post-edited sentence corrected by the corrector in the WMT’20 APE En-De and En-Zh tasks.

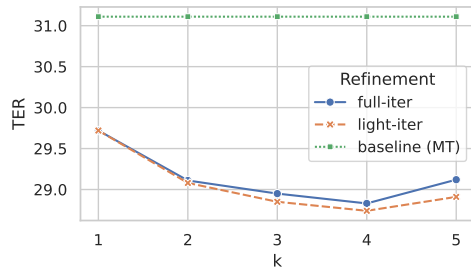
compared the translation quality in the first iteration. Table 6 shows the experimental results of detector-corrector with and without reordering. In TER, which indicates the number of edits to the reference translation, detector-corrector without reordering resulted in the same score as detector-corrector with reordering in En-De and degraded in En-Zh.

To investigate this gap in TER scores, we counted the total number of spans tagged by the detector and evaluated the TER score that measured the number of edits from the MT sentence to the post-edited sentence corrected by our detector-corrector (TER_{MT}). Table 7 shows that the number of edited spans was decreased by reordering, especially in En-Zh. In addition, the reordering operation reduces the TER_{MT} by 0.9% and 6.4% in En-De and En-Zh, respectively. This means that the number of edits from the MT sentence and the number of edits to the reference translation decreases by using the reordering operation; hence, the editing process becomes easier for humans to interpret.

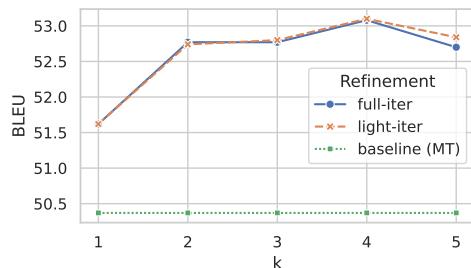
In summary, we confirmed that reordering is effective in reducing the number of edits, as shown by the TER scores in Table 6 and Table 7.

5.4 Effectiveness of Iterative Refinement

To verify the effectiveness of iterative refinement, we evaluated BLEU and TER scores in the WMT’20 En-De APE task at various numbers of inference iterations $k \in \{1, 2, 3, 4, 5\}$ on the development set. We also compared the difference between including (“full-iter”) and not including



(a) Comparison of TER scores for each iteration.



(b) Comparison of BLEU scores for each iteration.

Figure 4: Comparison of various iterations in iterative refinement. The scores were evaluated on the development set in the WMT’20 En-De APE task.

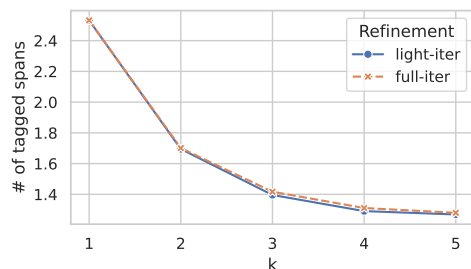


Figure 5: Number of tagged spans per sentence in the WMT’20 En-De APE task.

(“light-iter”) reordering when $k \geq 2$. Figure 4(a) and 4(b) shows that the first iterative refinement ($k = 2$) significantly improved the TER and BLEU scores from the first inference ($k = 1$). From $k = 2$ to 4, we see a slight improvement in both TER and BLEU. Comparing the iterative refinement methods, light-iter was slightly more accurate than full-iter, but the difference is lower than 0.1 % in both metrics.

Figure 5 shows the average number of bad- and insertion-tagged spans of MT sentences, which was corrected by the corrector. The figure shows that the number of corrected spans decreases in each iteration, especially when it significantly decreases in the second refinement, i.e., $k = 2$, which corresponds to the decrease of TER and BLEU in Figure 5.

Source	Georgia Lee , 89 , Australian jazz and blues singer .
Reference	乔治亚·李 (Georgia Lee) , 89 岁 , 澳大利亚 爵士 和 蓝调 歌手 。
MT (TER=64.7)	89 岁的 佐治亚州 李 , 澳大利亚 爵士乐 和 布鲁斯 歌手 。
Reordered MT	的 佐治亚州 李 89 岁 , 澳大利亚 爵士乐 和 布鲁斯 歌手 。
$k = 1$	
Annotated source	Georgia Lee <bad>, </bad> 89 , Australian jazz and blues singer .
Annotated MT	<bad>的</bad> 佐治亚 <bad>州</bad> 李 <ins></ins> 89 岁 , 澳大利亚 爵士乐和 <bad>布鲁斯</bad> 歌手 <bad>.</bad>
Correction	<bad></bad> <bad>·</bad> <ins>,</ins> <bad>蓝调</bad> <bad>。</bad>
Output (TER=35.3)	佐治亚·李 , 89 岁 , 澳大利亚 爵士乐 和 蓝调 歌手 。
$k = 2$	
Annotated source	Georgia Lee , 89 , Australian jazz and blues singer .
Annotated MT	佐治亚·李 <ins></ins> , 89 岁 , 澳大利亚爵士乐和蓝调歌手。
Correction	<ins> (George Lee) </ins>
Output (TER=17.7)	佐治亚·李 (George Lee) , 89 岁 , 澳大利亚 爵士乐 和 蓝调 歌手 。

Table 8: An example of the editing process.

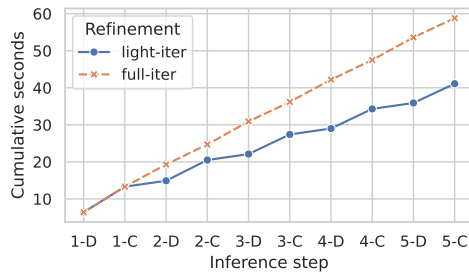


Figure 6: Cumulative time taken for each inference step. “ k -D” and “ k -C” denote the k -th inference step of the detector model and corrector model, respectively.

We also measured the cumulative time for each inference step. Figure 6 shows the total inference time in seconds for full-iter and light-iter when processing 1,000 sentences. In the figure, “ k -D” and “ k -C” denote the k -th inference step of the detector model and corrector model, respectively. It can be seen that light-iter infers faster than full-iter because light-iter does not predict reordering, which is time-consuming, in the detector inference at each iteration in $k \geq 2$.

From the results, our detector–corrector is further improved by using iterative refinement at least twice, and the inference speed is reduced by two-thirds using our lightweight iterative refinement without losing qualities.

5.5 Case Study: Editing Process

We analyzed examples of the editing processes of detector–corrector. Table 8 shows an example of the editing process of an MT sentence. In the table, the “Annotated source” line is the source sentences annotated with SRC-tag by the detector, and

the “Annotated MT” line is the reordered MT sentences annotated with MT-tag and MT-gap by the detector. The “Correction” and “Output” lines are the correction sequence generated by the corrector and the outputs of the detector–corrector, respectively. The table shows that our model detects and corrects the erroneous spans iteratively, and outputs the sentence with 17.7 TER in the second iteration. Note that the detector did not detect any erroneous spans in this example when $k \geq 3$. The table also shows that our model swaps two spans, “89 岁” and “佐治亚州 李”, which makes the word order align with the source sentence and reference translation.

6 Conclusion

We proposed “detector–corrector”, the edit-based automatic post-editing (APE) model, which explains which words are wrong in MT sentences and how to correct them for human post-editors. Experiments on the WMT’20 English–German and English–Chinese APE tasks showed that our detector–corrector model provides the editing process and outperformed the previous edit-based model, Levenshtein Transformer, and a black-box sequence-to-sequence APE model in TER.

In the future, we will further investigate what is needed to reduce the workload of human post-editors.

Acknowledgements

This work was partially supported by JSPS KAKENHI Grant Number JP21H05054.

References

- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In Bengio, Yoshua and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Bhattacharyya, Pushpak, Rajen Chatterjee, Markus Freitag, Diptesh Kanojia, Matteo Negri, and Marco Turchi. 2022. Findings of the WMT 2022 shared task on automatic post-editing. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 109–117, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Chatterjee, Rajen, Christian Federmann, Matteo Negri, and Marco Turchi. 2019. Findings of the WMT 2019 shared task on automatic post-editing. In Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 11–28, Florence, Italy, August. Association for Computational Linguistics.
- Chatterjee, Rajen, Markus Freitag, Matteo Negri, and Marco Turchi. 2020. Findings of the WMT 2020 shared task on automatic post-editing. In Barrault, Loïc, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Yvette Graham, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, and Matteo Negri, editors, *Proceedings of the Fifth Conference on Machine Translation*, pages 646–659, Online, November. Association for Computational Linguistics.
- Chen, Mengyun, Tao Ge, Xingxing Zhang, Furu Wei, and Ming Zhou. 2020. Improving the efficiency of grammatical error correction with erroneous span detection and correction. In Webber, Bonnie, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7162–7169, Online, November. Association for Computational Linguistics.
- Conneau, Alexis, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In Jurafsky, Dan, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online, July. Association for Computational Linguistics.
- Correia, Gonçalo M. and André F. T. Martins. 2019. A simple and effective approach to automatic post-editing with transfer learning. In Korhonen, Anna, David Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3050–3056, Florence, Italy, July. Association for Computational Linguistics.
- Cui, Qu, Shujian Huang, Jiahuan Li, Xiang Geng, Zaixiang Zheng, Guoping Huang, and Jiajun Chen. 2021. Directq: Direct pretraining for machine translation quality estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12719–12727.
- Deoghare, Sourabh and Pushpak Bhattacharyya. 2022. IIT Bombay’s WMT22 automatic post-editing shared task submission. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 682–688, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Deoghare, Sourabh, Diptesh Kanojia, Fred Blain, Tharindu Ranasinghe, and Pushpak Bhattacharyya. 2023. Quality estimation-assisted automatic post-editing. In Bouamor, Houda, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1686–1698, Singapore, December. Association for Computational Linguistics.
- Ding, Shuoyang, Marcin Junczys-Dowmunt, Matt Post, and Philipp Koehn. 2021. Levenshtein training for word-level quality estimation. In Moens, Marie-Francine, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6724–6733, Online and Punta

- Cana, Dominican Republic, November. Association for Computational Linguistics.
- Dou, Zi-Yi and Graham Neubig. 2021. Word alignment by fine-tuning embeddings on parallel corpora. In Merlo, Paola, Jorg Tiedemann, and Reut Tsarfaty, editors, *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2112–2128, Online, April. Association for Computational Linguistics.
- Dyer, Chris, Victor Chahuneau, and Noah A. Smith. 2013. A simple, fast, and effective reparameterization of IBM model 2. In Vanderwende, Lucy, Hal Daumé III, and Katrin Kirchoff, editors, *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648, Atlanta, Georgia, June. Association for Computational Linguistics.
- Fomicheva, Marina, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. 2020. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:539–555.
- Gu, Jiatao, Zhengdong Lu, Hang Li, and Victor O.K. Li. 2016. Incorporating copying mechanism in sequence-to-sequence learning. In Erk, Katrin and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1631–1640, Berlin, Germany, August. Association for Computational Linguistics.
- Gu, Jiatao, Changhan Wang, and Junbo Zhao. 2019. Levenshtein transformer. In Wallach, H., H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Herbig, Nico, Tim Düwel, Santanu Pal, Kalliopi Meladaki, Mahsa Monshizadeh, Antonio Krüger, and Josef van Genabith. 2020. MMPE: A Multi-Modal Interface for Post-Editing Machine Translation. In Jurafsky, Dan, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1691–1702, Online, July. Association for Computational Linguistics.
- Huang, Xuancheng, Yang Liu, Huanbo Luan, Jingfang Xu, and Maosong Sun. 2019. Learning to copy for automatic post-editing. In Inui, Kentaro, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6122–6132, Hong Kong, China, November. Association for Computational Linguistics.
- Huang, Xiaoying, Xingrui Lou, Fan Zhang, and Tu Mei. 2022. LUL’s WMT22 automatic post-editing shared task submission. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 689–693, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Jalili Sabet, Masoud, Philipp Dufter, François Yvon, and Hinrich Schütze. 2020. SimAlign: High quality word alignments without parallel training data using static and contextualized embeddings. In Cohn, Trevor, Yulan He, and Yang Liu, editors, *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1627–1643, Online, November. Association for Computational Linguistics.
- Junczys-Dowmunt, Marcin and Roman Grundkiewicz. 2018. MS-UEdin submission to the WMT2018 APE shared task: Dual-source transformer for automatic post-editing. In Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 822–826, Belgium, Brussels, October. Association for Computational Linguistics.
- Kasai, Jungo, James Cross, Marjan Ghazvininejad, and Jiatao Gu. 2020. Non-autoregressive machine translation with disentangled context transformer. In III, Hal Daumé and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 5144–5155. PMLR, 13–18 Jul.
- Kim, Hyun, Hun-Young Jung, Hongseok Kwon, Jong-Hyeok Lee, and Seung-Hoon Na. 2017a. Predictor-estimator: Neural quality estimation based on target word prediction for machine translation. *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*, 17(1), sep.
- Kim, Hyun, Jong-Hyeok Lee, and Seung-Hoon Na. 2017b. Predictor-estimator using multilevel task learning with stack propagation for neural quality estimation. In Bojar, Ondřej, Christian Buck, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, and Julia Kreutzer, editors, *Proceedings of the Second Conference on Machine*

- Translation*, pages 562–568, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Lee, Jihyung, WonKee Lee, Jaehun Shin, Baikjin Jung, Young-Kil Kim, and Jong-Hyeok Lee. 2020. POSTECH-ETRI’s submission to the WMT2020 APE shared task: Automatic post-editing with cross-lingual language model. In Barrault, Loïc, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Yvette Graham, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, and Matteo Negri, editors, *Proceedings of the Fifth Conference on Machine Translation*, pages 777–782, Online, November. Association for Computational Linguistics.
- Luong, Thang, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. In Márquez, Lluís, Chris Callison-Burch, and Jian Su, editors, *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1412–1421, Lisbon, Portugal, September. Association for Computational Linguistics.
- Mallinson, Jonathan, Aliaksei Severyn, Eric Malmi, and Guillermo Garrido. 2020. FELIX: Flexible text editing through tagging and insertion. In Cohn, Trevor, Yulan He, and Yang Liu, editors, *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1244–1255, Online, November. Association for Computational Linguistics.
- Mallinson, Jonathan, Jakub Adamek, Eric Malmi, and Aliaksei Severyn. 2022. EdiT5: Semi-autoregressive text editing with t5 warm-start. In Goldberg, Yoav, Zornitsa Kozareva, and Yue Zhang, editors, *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2126–2138, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.
- Malmi, Eric, Sebastian Krause, Sascha Rothe, Daniil Mirylenka, and Aliaksei Severyn. 2019. Encode, tag, realize: High-precision text editing. In Inui, Kentaro, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5054–5065, Hong Kong, China, November. Association for Computational Linguistics.
- Matthews, Brian W. 1975. Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405(2):442–451.
- Negri, Matteo, Marco Turchi, Rajen Chatterjee, and Nicola Bertoldi. 2018. ESCAPE: a large-scale synthetic corpus for automatic post-editing. In Calzolari, Nicoletta, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Koiti Hasida, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asuncion Moreno, Jan Odijk, Stelios Piperidis, and Takenobu Tokunaga, editors, *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan, May. European Language Resources Association (ELRA).
- Omelianchuk, Kostiantyn, Vitaliy Atrasevych, Artem Chernodub, and Oleksandr Skurzshanskiy. 2020. GECToR – grammatical error correction: Tag, not rewrite. In Burstein, Jill, Ekaterina Kochmar, Claudia Leacock, Nitin Madnani, Ildikó Pilán, Helen Yannakoudakis, and Torsten Zesch, editors, *Proceedings of the Fifteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 163–170, Seattle, WA, USA → Online, July. Association for Computational Linguistics.
- Ott, Myle, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation. In Dy, Jennifer G. and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 3953–3962. PMLR.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névéal, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October. Association for Computational Linguistics.
- Raffel, Colin, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- Ranasinghe, Tharindu, Constantin Orasan, and Ruslan Mitkov. 2020. TransQuest: Translation quality estimation with cross-lingual transformers. In Scott, Donia, Nuria Bel, and Chengqing Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5070–5081, Barcelona, Spain (Online), December. International Committee on Computational Linguistics.
- Ranasinghe, Tharindu, Constantin Orasan, and Ruslan Mitkov. 2021. An exploratory analysis of multilingual word-level quality estimation with cross-lingual transformers. In Zong, Chengqing, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the*

- 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 434–440, Online, August. Association for Computational Linguistics.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In Webber, Bonnie, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.
- Rei, Ricardo, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Nèveól, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Sharma, Abhishek, Prabhakar Gupta, and Anil Nelakanti. 2021. Adapting neural machine translation for automatic post-editing. In Barrault, Loïc, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussa, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Tom Kocmi, Andre Martins, Makoto Morishita, and Christof Monz, editors, *Proceedings of the Sixth Conference on Machine Translation*, pages 315–319, Online, November. Association for Computational Linguistics.
- Snover, Matthew, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA, August 8-12. Association for Machine Translation in the Americas.
- Specia, Lucia, Frédéric Blain, Marina Fomicheva, Erick Fonseca, Vishrav Chaudhary, Francisco Guzmán, and André F. T. Martins. 2020. Findings of the WMT 2020 shared task on quality estimation. In Barrault, Loïc, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Yvette Graham, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, and Matteo Negri, editors, *Proceedings of the Fifth Conference on Machine Translation*, pages 743–764, Online, November. Association for Computational Linguistics.
- Stahlberg, Felix and Shankar Kumar. 2020. Seq2Edits: Sequence transduction using span-level edit operations. In Webber, Bonnie, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5147–5159, Online, November. Association for Computational Linguistics.
- Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2, NIPS'14*, page 3104–3112, Cambridge, MA, USA. MIT Press.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Guyon, I, U V Luxburg, S Bengio, H Wallach, R Fergus, S Vishwanathan, and R Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.
- Vinyals, Oriol, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. In Cortes, C., N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Vu, Thuy-Trang and Gholamreza Haffari. 2018. Automatic post-editing of machine translation: A neural programmer-interpreter approach. In Riloff, Ellen, David Chiang, Julia Hockenmaier, and Jun'ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3048–3053, Brussels, Belgium, October-November. Association for Computational Linguistics.
- Wang, Jiayi, Ke Wang, Kai Fan, Yuqi Zhang, Jun Lu, Xin Ge, Yangbin Shi, and Yu Zhao. 2020. Alibaba's submission for the WMT 2020 APE shared task: Improving automatic post-editing with pre-trained conditional cross-lingual BERT. In Barrault, Loïc, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Yvette Graham, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, and Matteo Negri, editors, *Proceedings of the Fifth Conference on Machine Translation*, pages 789–796, Online, November. Association for Computational Linguistics.

Wu, Yonghui, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144.

Yang, Hao, Minghan Wang, Daimeng Wei, Hengchao Shang, Jiaxin Guo, Zongyao Li, Lizhi Lei, Ying Qin, Shimin Tao, Shiliang Sun, and Yimeng Chen. 2020. HW-TSC’s participation at WMT 2020 automatic post editing shared task. In Barrault, Loïc, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Yvette Graham, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, and Matteo Negri, editors, *Proceedings of the Fifth Conference on Machine Translation*, pages 797–802, Online, November. Association for Computational Linguistics.

Yang, Zhen, Fandong Meng, Yingxue Zhang, Ernan Li, and Jie Zhou. 2022a. Findings of the WMT 2022 shared task on translation suggestion. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 821–829, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.

Yang, Zhen, Fandong Meng, Yingxue Zhang, Ernan Li, and Jie Zhou. 2022b. WeTS: A benchmark for translation suggestion. In Goldberg, Yoav, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5278–5290, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.

A Ethical Considerations

We trained all models from open datasets; therefore, if their datasets have toxic text, the models may have the risk of generating toxic content.

B Limitations

Our model can show the editing process and correction candidates by taking into account the opinions of professional translators, but we have not conducted a human evaluation of how much they affect the actual post-editing process.

Our method may demand a larger memory footprint than a single seq2seq model because it runs two models, the detector and corrector.

Our study focuses on correcting translation errors, and thus our model cannot detect and correct non-factual information when including them in a source sentence.

Our model only corrects the erroneous spans detected by the detector; thus, spans that the detector fails to detect may remain uncorrected.

C Tools, Models, and Datasets

Tools We implemented all models in FAIRSEQ which is published under the MIT-license.

Models We used the following pre-trained NMT models implemented in FAIRSEQ to create the training data.

- En-De: https://www.quest.dcs.shef.ac.uk/wmt20_files_qe/models_en-de.tar.gz
- En-Zh: https://www.quest.dcs.shef.ac.uk/wmt20_files_qe/models_en-zh.tar.gz

Our models were trained by using NVIDIA A6000 GPU. The training costs, “GPU hours”, multiplied by the number of GPUs and computation time, are shown in Table 9. Note that the translation performance for each model was evaluated with only a single training.

Datasets We evaluated all models using WMT’20 APE datasets published under the Creative Commons Zero v1.0 Universal license. Parallel data of the WMT’19 En-De and En-Zh translation tasks, used in our training data, can be used for research purposes as described in <https://www.statmt.org/wmt19/translation-task.html>.

In the En-Zh task, we tokenized the test set of the En-Zh APE task using JIEBA⁷ to calculate the TER and BLEU scores.

⁷<https://github.com/fxsjy/jieba>

Seq2Seq	
Encoder	XLM-R large (24 layers)
Decoder	Transformer decoder
Number of layers	6
Hidden size	1024
FFN hidden size	4096
Learning rate	1e-4
Batch size	24,000 tokens
Training steps	60,000
Training cost	24.6 GPU hours
LevT	
Encoder	XLM-R large (24 layers)
Decoder	Transformer decoder
Number of layers	6
Hidden size	1024
FFN hidden size	4096
Learning rate	1e-4
Batch size	12,000 tokens
Training steps	60,000
Training cost	12.4 GPU hours
Detector	
Encoder	XLM-R large (24 layers)
Decoder	Transformer decoder
Number of layers	4
Hidden size	1024
FFN hidden size	4096
Learning rate	3e-5
Batch size	6,000 tokens
Training steps	40,000
Training cost	8.0 GPU hours
Corrector	
Encoder	XLM-R large (24 layers)
Decoder	Transformer decoder
Number of layers	6
Hidden size	1024
FFN hidden size	4096
Learning rate	1e-4
Batch size	24,000 tokens
Training steps	60,000
Training cost	29.0 GPU hours

Table 9: Hyperparameters of the models.

The statistics of the training data are shown in Table 10.

	DAug for detector	
	w/o	w/
(1) APE task data	7,000	7,000
(2) Translation task data	2,000,000	2,000,000
<i>Training data of detector</i>		
Base data: (1)×20 + (2)	2,140,000	4,280,000
<i>Training data of corrector</i>		
Base data: (1)×20 + (2)	2,140,000	4,280,000
+ MT training	4,280,000	8,560,000
+ PE training	4,280,000	8,560,000
+ MT & PE training	6,420,000	12,840,000

Table 10: Statistics of the training data. “DAug” denotes data augmentation. In the experiment, to make the difference in data size fair, we trained with the same number of parameter updates without using the number of epochs, i.e., the number of training epochs decreases as the data size increases.

Assessing Translation Capabilities of Large Language Models involving English and Indian Languages

Vandan Mujadia, Ashok Urlana, Yash Bhaskar, Penumalla Aditya Pavani,
Kukkapalli Shravya, Parameswari Krishnamurthy, Dipti Misra Sharma

IIT Hyderabad, India

{vandan.mu, ashok.urlana, yash.bhaskar}@research.iiit.ac.in,

{aditya.pavani, kukkapalli.shravya}@students.iiit.ac.in,

{param.krishna, dipti}@iiit.ac.in

Abstract

Generative Large Language Models (LLMs) have achieved remarkable advances in various NLP tasks. In this work, our aim is to explore the multilingual capabilities of large language models by using machine translation as a task involving English and 22 Indian languages. We first investigate the translation capabilities of raw large-language models, followed by exploring the in-context learning capabilities of the same raw models. We fine-tune these large language models using parameter-efficient fine-tuning methods such as LoRA and additionally with full fine-tuning. Through our study, we have identified the model that performs best among the large language models available for the translation task.

Our results demonstrate significant progress, with average BLEU scores of 13.42, 15.93, 12.13, 12.30, and 12.07, as well as chrF scores of 43.98, 46.99, 42.55, 42.42, and 45.39, respectively, using two-stage fine-tuned LLaMA-13b for English to Indian languages on IN22 (conversational), IN22 (general), flores200-dev, flores200-devtest, and newstest2019 testsets. Similarly, for Indian languages to English, we achieved average BLEU scores of 14.03, 16.65, 16.17, 15.35 and 12.55 along with chrF scores of 36.71, 40.44, 40.26, 39.51, and 36.20, respectively, using fine-tuned LLaMA-13b on IN22

(conversational), IN22 (general), flores200-dev, flores200-devtest and newstest2019 testsets. Overall, our findings highlight the potential and strength of large language models for machine translation capabilities, including languages that are currently underrepresented in LLMs.

1 Introduction

Generative Large Language Models (LLMs) have made significant performance improvements in various natural language processing (NLP) tasks, showcasing exceptional progress in a wide range of applications (Xuanfan and Piji, 2023; Xi et al., 2023). These tasks range from open domain question answering, where LLMs excel in providing accurate and coherent responses, to instruction-based tasks such as code completion, where LLMs can generate code snippets based on given prompts (Vaithilingam et al., 2022). LLMs have also demonstrated proficiency in tasks such as writing essays, grammar checks (Wu et al., 2023a), and text summarization, where they can produce high quality results (Chang et al., 2023). These advances have been observed mainly in tasks centered on English. The popular LLMs support several natural languages. The performance of some languages other than English is not yet on par or yet to be evaluated (Lai et al., 2023; Zhu et al., 2023).

A multilingual country like India, where over 364+ languages and dialects¹ are spoken across its vast territory, presents a multitude of challenges across various domains due to language barriers (Zieliński and others, 2021), such as day-to-day communication, education (Steigerwald et al., 2022), business, healthcare (Mehandru et al., 2022),

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹https://en.wikipedia.org/wiki/Linguistic_Survey_of_India

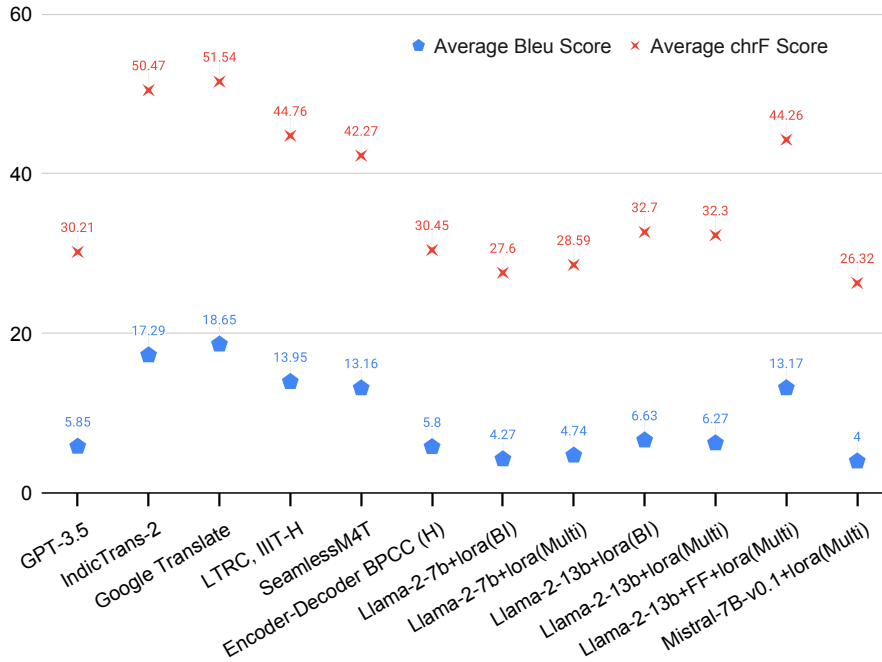


Figure 1: LLMs based Machine Translation performance comparison with public systems for **English to Indian Languages**. BLEU and chrF scores are averaged over 22 Indian Languages and 5 different benchmark data sets. The available MT systems are GPT-3.5 (GPT-3.5 Davinci, by OpenAI), IndicTrans-2, Google Translate, LTRC-IIIT-H, SeamlessM4T. LLaMA-2-7b and LLaMA-2-13b are evaluated as LLM based fine-tuned MT systems are namely LLaMA-2-7b+lora (Multi), LLaMA-2-13b+lora (Multi), and LLaMA-2-13b+FF+lora (Multi). Encoder-Decoder BPCC (H) represents scores for encoder decoder based transformer model trained on BPCC Human Training data.

tourism, governance, and more. Recent advancements in the field of Large Language Models may offer solutions to these challenges tailored to Indian languages.

To test whether decoder-based LLMs can effectively overcome language barriers, it is crucial to assess the proficiency of large language models in handling Indian languages. Machine Translation, as a critical multilingual task, could be an ideal option to explore the multilingual capabilities of existing models. Hence, we can formulate the question to assess the proficiency of large language models in handling Indian languages as follows: **How effectively do large language models perform in multilingual tasks like Machine Translation, particularly when dealing with Indian languages?**

In this work, our main contribution is to address the following points in response to the above question.

- What are the directions for utilizing or adapting available Large Language Models for Indian Languages?
 - How do LLMs perform in translating a wide range of Indian languages un-

der zero-shot and in-context learning settings?

- Does LLM fine-tuning improve the translation capabilities of Large Language Models? How do they perform in low-resourced MT languages?
- The Impact of LLM Vocabulary on the Performance of Large Language Models in Translation Tasks.

To address the above questions, we assess the translation capabilities of popular large language models (opt, bloom, LLaMA-1, MPT, Falcon, LLaMA-2, and Mistral (§B)) that involve English and 22 scheduled Indian languages (Assamese, Bangla, Bodo, Dogri, Konkani, Gujarati, Hindi, Kannada, Kashmiri, Maithili, Malayalam, Marathi, Meitei, Nepali, Odia, Punjabi, Sanskrit, Santali, Sindhi, Tamil, Telugu, and Urdu). We initially examine the translation capabilities of these raw large language models mentioned above (§4.1). Subsequently, we explore their in-context learning abilities (§4.1). Additionally, we fine-tune the base models using parameter-efficient fine-tuning methods, specifically LoRa (§5). Furthermore, we investigate

the potential of two-stage fine-tuning for large language models, which involves full parameter fine-tuning in the first stage, followed by LoRa-based adapter fine-tuning (§5).

The key findings of our work, summarized in Figure 1, highlight the performance of our LLM-based machine translation fine-tuned models compared to various known translation engines. These engines range from commercial (Google², GPT-3.5³) to open source (IndicTrans-2⁴, LTRC-IIIT-H⁵, seamless4t⁶), our trained Encoder-Decoder BPPC (H) model (Appendix A), traditional supervised encoder-decoder translation models (Google, IndicTrans-2, LTRC-IIIT-H), and decoder-driven causal large language model-based translation systems (GPT-3.5).

Our findings underscore the significant potential of large language models for translation tasks involving English and Indian Languages. Although raw LLMs (LLaMA-2-7b and LLaMA-2-13b) do not perform well in translation tasks, our two-stage MT fine-tuned models (LLaMA-2-13b + FF + lora (Multi)) yield comparative results even with minimal parallel corpora. This suggests that LLMs have the potential to possess multilingual capabilities for translating into underrepresented languages, which can be further enhanced by fine-tuning. This work will be a crucial and pioneering milestone in evaluating LLMs for language representation and assessing their translation capabilities for a diverse range of Indian languages, especially those with limited available resources.

2 Related Work

Recent advancements in machine translation have shown that neural machine translation (NMT) has made significant strides in terms of output fluency and translation quality, especially when ample parallel data are available (Barrault et al., 2020). However, the scarcity or absence of parallel data poses a challenge for most language pairs. In the case of Indian languages, recent developments have tried to address this issue by introducing a new state-of-the-art approach: multilingual machine translation involving Indian languages and English (Wang et al., 2021; Dabre et al., 2020; Madaan and Sadat,

2020). This approach uses a single script for machine translation, taking advantage of the lexical and syntactic similarities that arise from genetic and contact relatedness among Indian languages (Gala et al., 2023; Eriguchi et al., 2022; Bapna and Firat, 2019).

In the field of LLM driven machine translation, in-context learning has gained significant attention (Wu et al., 2023b). The use of large language models (LLMs) for multilingual machine translation has been a topic of interest (Zhang et al., 2023). Recent studies have evaluated the translation capabilities of LLMs for different language directions, with a focus on models like ChatGPT (Bang et al., 2023).

In particular, (Xu et al., 2023) proposed a two-stage fine-tuning approach for machine translation using LLM, involving fine-tuning in monolingual data followed by fine-tuning on a small set of high-quality parallel data. Our work represents the first study to specifically explore machine translation involving Indian languages using large language models. The details on Large Language Models are presented in the Appendix B.

3 Indian Languages representation in LLMs

Pre-trained (or Base/Raw) large language models are trained on a huge amount of language data, and some of these models are trained on multiple languages (Naveed et al., 2023). However, their training mainly focuses on the English text (Penedo et al., 2023a). The emphasis on English is due to its substantial presence on the Internet and its widespread usage in business contexts. For the purpose of this work, our objective is to assess the effectiveness of these models in Machine Translation tasks that involve both English and Indian Languages. Consequently, it becomes crucial to investigate the representation of Indian languages within these large language models.

An approach to investigating the representation of Indian languages within a large language model can involve analyzing the frequency of language-specific words and sentences used during the training of these models. Unfortunately, it is not possible to perform this analysis as the training data used for these models are not publicly accessible. LLaMA-2, in particular, has mentioned that its pretraining corpus consists mainly of English and may not be optimal for other languages (Touvron et al., 2023b). However, it is worth mentioning that approximately

²<https://translate.google.co.in/>

³<https://chat.openai.com/>

⁴<https://github.com/AI4Bharat/IndicTrans2>

⁵<https://ssmt.iit.ac.in/translate>

⁶https://github.com/facebookresearch/seamless_communication

Language Family Language	Indo-Aryan										Dravidian				Sino-Tibetan		Austroasiatic								
	asm	ban	kas	snd	urd	doi	hin	gom	mai	mar	nep	san	guj	odi	pan	kan	mal	tam	tel	mini	brx	sat			
Language Script	Bangla		Perso-Arabic			Devanagari										Kannada	Malayalam	Tamil	Telugu	Meitei	Devanagari	Ol Chik			
No of Letters in Unicode	96		256			128										91	91	80	91	118	72	100	56	96	48
Models (Vocab)																									
BLOOM (250680)	(48,48)	(49,207)	(67,61)	(57,34)	(56,35)	(55,25)	(62,29)	(66,52)	(46,26)	(61,39)	(00,56)	(67,29)	(00,48)												
FALCON (65024)	(00,96)	(12,244)	(2,126)	(00,91)	(00,91)	(00,72)	(0,100)	(00,56)	(02,70)	(04,96)	(00,56)	(02,94)	(00,48)												
LLAMA-1.2 (32024)	(24,72)	(45,211)	(38,90)	(01,90)	(00,91)	(04,76)	(02,89)	(33,155)	(19,53)	(01,99)	(00,56)	(38,58)	(00,48)												
MISTRAL (32052)	(34,62)	(47,209)	(43,85)	(05,86)	(00,91)	(02,78)	(18,73)	(04,116)	(22,50)	(11,89)	(00,56)	(43,53)	(00,48)												
MPT (50277)	(05,91)	(35,221)	(22,106)	(02,89)	(00,91)	(00,80)	(00,91)	(01,117)	(05,67)	(03,97)	(00,56)	(22,74)	(00,48)												
OPT (50265)	(00,96)	(13,243)	(1,127)	(00,91)	(00,91)	(00,80)	(00,91)	(0,118)	(00,72)	(0,100)	(00,56)	(01,95)	(00,48)												

Table 1: The language support of various LLMs for 22 Indian languages, along with the corresponding families, scripts, and letters representing each language. In each tuple (xx, yy), the first value, xx represents the number of language-specific characters present in respective LLM, while the second value, yy indicates the number of language-specific characters supported in the form of bytes in respective LLM and for the respective language.

8.38% of the data includes languages other than English and codes in LLaMA-2.

On the other hand, studying the vocabulary (or letters/characters) of a corpus can provide valuable insights into the representation and coverage of language within that corpus. The writing system or script used plays a crucial role in representing a language. Therefore, the analysis of the vocabulary can be considered a proximal task. Fortunately, we have access to the sub-word vocabulary for the considered large language models. By comparing the characters present in the subword vocabulary with those in the corresponding language script, we can approximate the language representation within the respective LLM.

For this work, we include a total of 22 scheduled Indian languages for translation, which can be categorized into four main language families: Indo-Aryan, Dravidian, Sino-Tibetan, and Austroasiatic. These 22 Indian languages are written using 13 major scripts. It is interesting to note that most of these scripts can be traced back to the Brahmi script⁷, which served as the basis for the development of several Indian scripts (Salomon, 1995). Each of these 13 writing systems has its own unique set of letters and characters⁸, reflecting the phonetic and linguistic characteristics of the respective languages they represent.

Table 1 presents an overview of the scripts, the languages that use these scripts, and the corresponding vocabulary sizes of the subwords for LLMs. The numbers indicated in ‘(X,Y)’ represent the counts of native script letters (characters in unicode⁹) present and not present in the respective LLM. Specifically, X denotes the number of characters in the native language that are present in the vocabulary, while Y denotes the number of characters

represented as predefined (multiple) hexadecimal values. In other words, there is no direct representation for these many Unicode characters. On analysis, we observe that, in general, the 22 Indian languages have a very limited presence in most LLMs. However, the Devanagari, Perso-Arabic, and Bangla scripts demonstrate a few subword vocabularies among 22 Indian Languages, whereas other scripts have minimal or near-zero representation within the vocabulary.

4 Experiment setup: Machine Translation using LLMs

To evaluate the performance of large language models (LLMs) in machine translation tasks involving English and 22 Indian languages, we mainly conducted two experiments. The first experiment focused on assessing the performance of the pre-trained (raw) LLM, and example-based in-context learning for the machine translation. In the second experiment, we explore the fine-tuning of the best-performing large language models for the translation task. Both directions of translation were explored, including English to 22 Indian languages and 22 Indian languages to English. All experiments were carried out using translation benchmark data, as discussed in Section 5.

As part of our experimental setup, we use the prompt pipeline shown in Figure 2. This pipeline involved using a Prompt Generator to generate specific prompts for the source and target language along with the source text. Subsequently, an LLM call is triggered to generate a response, which was then processed by a translation parser to obtain the actual translation. To ensure high-throughput and memory-efficient inference and serving of LLM, we utilized the vLLM library¹⁰ (Kwon et al., 2023). We conducted all experiments using a temperature parameter of 0, which ensures that the model be-

⁷<https://tinyurl.com/2r4zjd2d>

⁸https://en.wikipedia.org/wiki/Official_scripts_of_the_Republic_of_India

⁹<https://unicode.org/>

¹⁰<https://github.com/vllm-project/vllm>

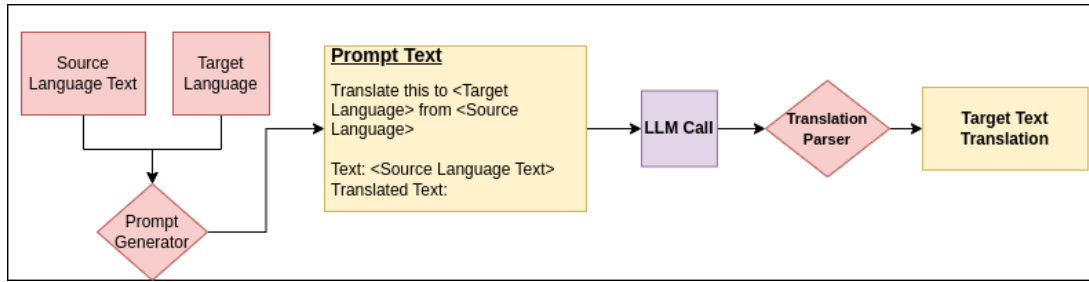


Figure 2: Prompting Mechanism for Translation

has deterministically. When setting the temperature to 0, the model is restricted to selecting the word with the highest probability, effectively limiting its choice to the most likely option (Aksitov et al., 2023). All of our experiments are conducted using the vLLM library on A100, 40GB GPUs.

4.1 Machine Translation on Raw LLM

To optimize the machine translation task on our selected LLMs, we performed manual trials with various prompts. Through these trials, we found that directly asking for translation and presenting the text in JSON format yielded better results, as the models seemed to comprehend the JSON structure more effectively (Reinauer et al., 2023). After multiple iterations, we finalized two prompts to translate sentences using raw (pretrained) LLMs, as illustrated in the following examples. These prompts were used to evaluate the efficiency of the models.

Example: Translation Prompt-1

Translate this to <Target Language> from <Source Language>

Text: <Source Language Text>
Translated Text:

Similarly, we identified and modified the prompt for example-based in-context learning with LLM. This prompt is specified in Example above (ICL Translation Prompt). In the case of in-context learning, all of our experiments involved providing a single translation sample as a contextual learning example prior to the actual translation command. We ensured that this example remained consistent for the same language pair in all LLM calls. The sample itself was randomly selected from the Human-BPCC translation training corpus (AI4Bharat et al., 2023). We present the results of both of these experiments in the Performance and Discussions

section.

Example: Translation Prompt-2

Translate this from <Source Language> to <Target Language>

<Source Language>: <Source Language Text>
<Target Language>:

4.2 Fine-tuning LLM for Machine Translation

To examine the potential improvement in multilingual understanding or translation performance of LLM beyond the baseline pre-trained LLM (Raw / Base), we conducted fine-tuning experiments for the translation task.

Example: ICL Translation Prompt

If the <Source Language> to <Target Language> translation for “<Source Example>” is “<Target Example>” from <Source Language>, following that, translate this to <Target Language> from <Source Language>

Text: <Source Language Text>
Translated Text:

4.2.1 Training Data

To fine-tune large language models (LLMs) for the machine translation task, we used the Bharat Parallel Corpus Collection (BPCC) (AI4Bharat et al., 2023). This corpus is publicly available and specifically translated in English to 22 Indian languages. It consists of two main parts: BPCC-Mined and BPCC-Human, comprising a total of approximately 230 million parallel text pairs. For the fine-tuning process, we focus on the BPCC-Human

English-	#Sents	S-AvgL	T-AvgL	S-Words	T-Words	S-Types	T-Types
<i>Assamese (asm)</i>	138208	16.88	14.39	2333583	1988395	125480	185151
<i>Bangla (ban)</i>	180219	17.80	15.07	3208203	2715959	161820	227468
<i>Bodo (brx)</i>	113139	17.79	13.96	2012274	1579042	116963	227180
<i>Dogri (doi)</i>	24157	15.32	17.68	370047	427110	48256	41370
<i>Konkani (gom)</i>	97555	17.13	14.03	1671465	1368512	82783	145300
<i>Gujarati (guj)</i>	135664	17.71	15.96	2402552	2164831	123935	174886
<i>Hindi (hin)</i>	222356	17.84	19.69	3966247	4378231	183737	202423
<i>Kannada (kan)</i>	117222	16.83	12.44	1972881	1458053	100778	208803
<i>Kashmiri (kas)</i>	19824	16.02	17.68	317634	350577	43197	66210
<i>Maithili (mai)</i>	23690	16.11	15.79	381720	374042	52920	57423
<i>Malayalam (mal)</i>	137950	16.30	11.13	2248081	1535654	120999	299146
<i>Marathi (mar)</i>	175893	17.94	14.81	3154904	2604119	167822	299983
<i>Meitei (mni)</i>	56617	17.77	15.73	1006271	890828	86175	161043
<i>Nepali (nep)</i>	85442	16.76	14.13	1431858	1207687	105411	145175
<i>Odia (odi)</i>	36923	17.07	15.49	630148	571958	68765	79932
<i>Punjabi (pan)</i>	80951	17.22	18.29	1394286	1480835	63510	74451
<i>Sanskrit (san)</i>	33189	16.30	11.69	541034	387957	61591	119856
<i>Santali (sat)</i>	24368	16.95	19.28	412918	469791	51307	56053
<i>Sindhi (sin)</i>	10503	17.10	19.32	179592	202952	28945	30782
<i>Tamil (tam)</i>	150254	17.76	13.34	2668252	2004981	139214	290917
<i>Telugu (tel)</i>	111808	16.81	12.64	1879737	1413466	96105	191792
<i>Urdu (urd)</i>	150747	17.62	20.20	2656814	3044480	144001	129856

Table 2: English to Indian Languages machine translation Fine-tuning data from BPC-Human (AI4Bharat et al., 2023). In this, the term “#Sents” refers to the total number of parallel sentences. “S-AvgL” and “T-AvgL” represent the average sentence length, in terms of words, for the source and target languages, respectively. Likewise, “Words” denotes the total number of words, while “Type” represents the total number of unique words.

Method	Hyper-parameter	Value
LoRA	LoRA modules	PEFT ¹¹
	rank	8
	dropout	0.05
	learning rate	1e-4
	global batch size	8
Full-parameter FSDP	epochs	6
	learning rate	1e-4
	global batch size	4
	epochs	5

Table 3: Hyper-parameter configurations of LoRA based and full fine-tuning for 4*A100 40GB GPUs

dataset, which contains 2.2 million English-Indic pairs. Additionally, this data set includes subsets derived from sentences from English Wikipedia and everyday usage scenarios. For more information on this corpus, we present Table 2. It shows a diverse representation of multilingual parallel corpora in terms of sentence length and the number of characters per token (compare T-AvgL with S-AvgL) for 22 Indian languages.

4.2.2 Fine-tuning Details

Considering the raw LLM performance, model parameters, and resource constraints, we selected a subset of LLMs for the fine-tuning process. Specifically, we chose LLaMA-2-7b, LLaMA-2-13b, and Mistral-7B for the fine-tuning experiment. For the selected LLMs, we decided to perform fine-tuning considering multiple options to check their performance. These options included bilingual translation fine-tuning, multilingual translation fine-tuning, low-rank adaptation-based fine-tuning, and a two-stage fine-tuning approach: full fine-tuning followed by low-rank adaptation-based fine-tuning. Due to limitations in training resources, we prioritize full fine-tuning only for best performing Large Language Models.

Specifically, we performed LoRa-based fine-tuning (Hu et al., 2021) for all English to 22 Indian languages (in both directions) in bilingual settings using LLaMA-2-7b and LLaMA-2-13b. Furthermore, we performed LoRa-based multilingual fine-tuning for English to the combined 22 Indian languages, as well as for the combined 22 Indian

TestSet	#Sent	Details
IN22_conv_test IN22_gen_test	1502 1023	(AI4Bharat et al., 2023) released MT benchmark data covering English to 22 Indian Languages.
Flores200-dev Flores200-devtest	997 1012	(Goyal et al., 2022) released MT benchmark data which includes English to 17 Indian Language pairs considered in this work.
Newstest2019	1997	(Federmann et al., 2022) released MT benchmark data which includes English to 10 Indian Language pairs considered in this work.

Table 4: Benchmark data details covering English to 22 Indian Languages

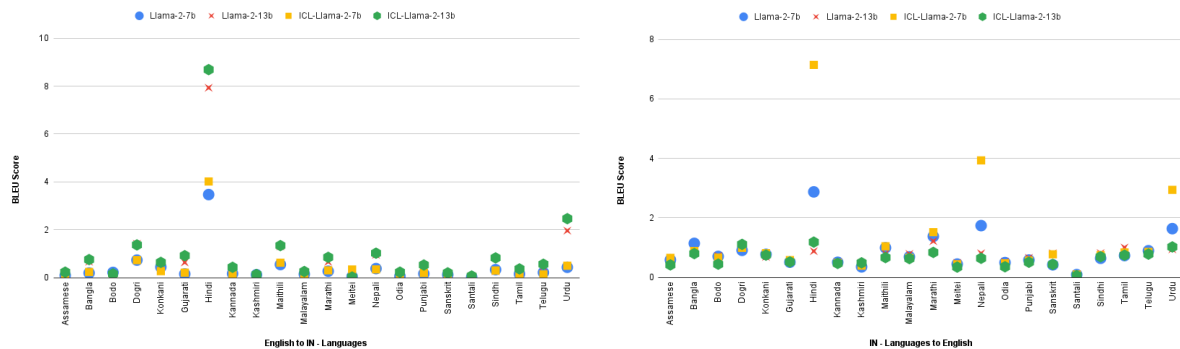


Figure 3: Evaluation of English - 22 Indic language Translation over 5 benchmark-sets (averaged): Raw LLM vs. In-Context Learning (ICL); Raw LLM models: LLaMA-2-7b, LLaMA-2-13b)

languages to English, using LLaMA-2-7b, LLaMA-2-13b and Mistral-7B. Based on the overall performance, we proceeded with a two-stage fine-tuning approach for the multilingual translation task specifically on LLaMA-2-13b. In the first stage, we performed a full fine-tuning for the multilingual translation objective. Subsequently, in the second stage, we performed LoRa-based fine-tuning based on same multilingual translation tasks on the fully fine-tuned model.

For both types of LLM fine-tuning, we utilize the llama-recipes codebase¹² that provides an efficient implementation for LoRa-based adapter fine-tuning with PEFT (Mangrulkar et al., 2022). For more details, the llama-recipes documentation can be referred to. Hyperparameters for the fine-tuning process are specified in Table 3. Training data used for fine-tuning experiments will be presented in Subsection 4.2.1.

5 Machine Translation Benchmark Data

We evaluated the performance of multilingual translation using three different benchmark datasets, as outlined in Table 4. The table provides a comprehensive overview of each translation benchmark dataset, highlighting the availability of n-way parallel data for the specified number of Indian languages from English as the source direction.

¹²<https://github.com/facebookresearch/llama-recipes/>

6 Performance Evaluation

We evaluated the performance of the translation outputs using the BLEU (Papineni et al., 2002) and chrF (Popović, 2015) evaluation methods on benchmark data described in Section 5. However, we did not include COMET (Rei et al., 2022) as an evaluation method due to the absence of support for many low-resource Indian languages at the time of evaluation. We used sacreBLEU library (Post, 2018) for BLEU¹³ and chrF¹⁴ calculation. To mitigate the impact of randomness in scores, we present our findings as the average of two runs for all our results.

Raw (Zero shot) vs ICL based Translation on LLMs Figure 3 presents the comparison of the overall results when evaluating the quality of translation for LLM outputs based on raw LLM and In Context Learning (ICL) based LLM outputs. The left subfigure represents the results for English to 22 Indian languages, while the right subfigure presents the results for 22 Indian languages to English translation. We observed amplified performance for the Bloom large language model for certain languages, which can be attributed to the known leak of MT benchmark data in the pretraining (Zhu et al., 2023).

¹³footprint for BLEU: nrefs:1—case:mixed—eff:no—tok:13a—smooth:exp—version:2.1.0

¹⁴footprint for chrF: nrefs:1—case:mixed—eff:yes—nc:6—nw:0—space:no—version:2.1.0

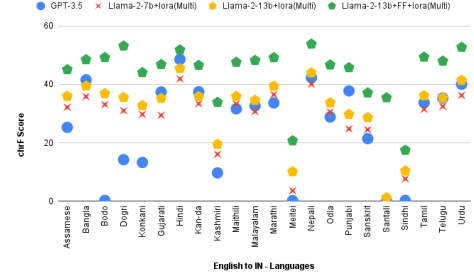
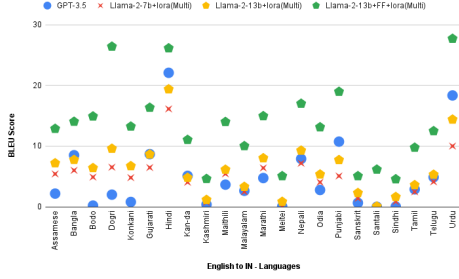


Figure 4: Performance comparison of GPT-3.5 vs our Fine-Tuned LLM Translation models (LLaMA-2-7b+lora (Multi), LLaMA-2-13b+lora (Multi), and LLaMA-2-13b+FF+lora (Multi)): English to 22 Indian languages over 5 benchmark sets (averaged). Here, LORA stands for Low-Rank Adaptation of Large Language Models-based fine-tuning. Multi stands for the multilingual model, FF stands for full fine tuning, and FF + lora stands for 2-stage fine tuning.

Consequently, that was the reason for excluding this language model from further experiments.

LLM models such as OPT, MPT, LLAMA-1 and Falcon exhibited poor performance, which can be correlated with the absence or minimal presence of characters for our focused Indian Languages in their vocabulary (Table 1). Therefore, we have omitted reporting the results for these models. Figure 3 indicates that Llama-2 models show relatively better performance with ICL settings compared to the raw models. Detailed results are presented in the appendix D.

Through manual analysis, we observed that less-represented (in vocabulary) languages such as Gujarati, Kannada, Odia, etc. (Table 1), ICL-driven translation tends to repeat the same translation given in the context of learning. On the other hand, raw models tend to hallucinate and repeat words throughout translation (Guerreiro et al., 2023) for these languages.

An important finding from manual analysis is that these raw LLMs demonstrate the ability to accurately identify languages (e.g., when asked for Gujarati translation, it gives inaccurate translations but correctly hallucinates text in the Gujarati script). This is a positive aspect and indicates a significant advantage of these LLMs in terms of their understanding and differentiation of languages and language scripts. In response to the question asked in the Introduction, it is true that the major LLMs available are primarily focused on English. However, *they exhibit minimal potential for zero shot and example-based translation capabilities.*

Fine-Tuned LLM driven Translations: English to Indian Languages We conducted an evaluation to compare the performance of our Fine-Tuned LLM models with GPT-3.5, as both models use

the same decoder-based approach in architecture. Figure 4 illustrates the comparison of English to 22 Indian language translations in terms of the BLEU and CHRF scores. The scores for GPT-3.5 are generally lower compared to our fine-tuned methods; also, our fine-tuned models have higher numbers than our previously mentioned zero-shot and example-based learning baseline. This indicates that with minimal parallel translation corpora, we are able to achieve considerable translations for translating into Indian languages from English.

Furthermore, we observed that multilingual fine-tuning yielded a better overall performance compared to bilingual fine-tuning. The two-stage fine-tuning approach also outperformed other fine-tuning methods for the translation task. The impressive results of the two-stage fine-tuning approach, as shown in Figures 4 and 1, are comparable to those of traditional encoder-decoder-based translation models. Note that this performance improvement was achieved using only a few thousand parallel data (Encoder-Decoder BPCC (H) model vs. LLM-based models in Figure 1), whereas traditional NMT models typically require a larger amount of data. From Figure 4, we can see that translating to low-resource languages such as Dogri, Konkani, Kashmiri, Meitei, Sanskrit and Sindhi yielded favorable evaluation numbers (Detailed results are presented in the Appendix D) compared to existing translation systems. In response to the question posed in the Introduction, the fine-tuning of LLM *enhances translation capabilities, particularly when using multilingual fine-tuning. These models also demonstrate proficiency in translating low-resource languages.*

Fine-Tuned LLM driven Translations: Indian Languages to English Figure 5 shows the com-

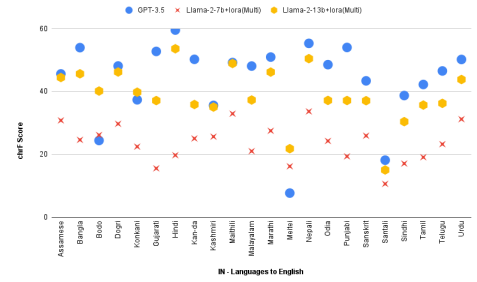
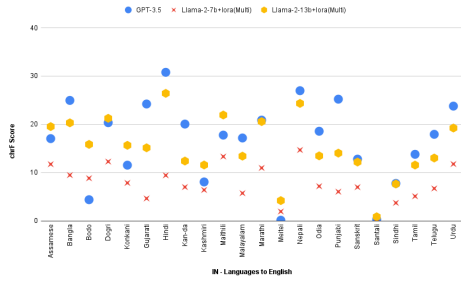


Figure 5: Performance comparison of GPT-3.5 vs our Fine-Tuned LLM Translation models (LLaMA-2-7b+lora (Multi), LLaMA-2-13b+lora (Multi), and LLaMA-2-13b+FF+lora (Multi)) on BLEU and CHRF scores: English to 22 Indian languages over 5 benchmark sets (averaged). Here, LORA stands for Low-Rank Adaptation of Large Language Models-based fine-tuning. Multi stands for the multilingual model.

parison of Indian language to English translation. The scores for GPT-3.5 are generally higher compared to our fine-tuned methods, while our fine-tuned models still outperform the previously mentioned zero-shot and example-based context learning-driven LLM results. In particular, the performance improvement for the Indian-language to English translation is comparatively lower than that for the English-to-Indian-language translation. Compared to translations from English to Indian languages, the LoRa-based single-stage fine-tuning here performs the best among all the fine-tuning approaches. Detailed results are presented in the appendix D.

This disparity can also be attributed to the representation of the Indian language vocabulary in these LLMs. As presented in Table 1, the subword vocabulary for Indian languages is limited in the LLMs considered. Consequently, when input is processed in Indian languages, characters that are not present in the vocabulary receive multiple hexadecimal representations of the vocabulary. This creates a bottleneck in finding the correct representation and, hence, the underlying meaning, making it challenging for the LLM network to perform the corresponding semantic translations. However, this issue is not prominent when translating from English to Indian languages, as the underlying understanding of English is robust for these large language models. This enables the network to effectively map the respective language translations.

6.1 Human Evaluation

For human machine translation evaluation, we used the direct assessment (DA) method (Stanchev et al., 2020). This method enables human evaluators to directly rate translations based on predetermined quality criteria. It involved a meticulous analysis

and comparison of machine-generated output with the source text, resulting in a continuous scale score ranging from 1 to 100. A score of 1 signifies non-sensical output, while a score of 100 indicates a perfect translation. This method provides a more objective and reliable assessment of the quality of machine translation.

For our evaluation, we conducted a direct assessment (Stanchev et al., 2020) for four language pairs: English to Hindi, Marathi, Tamil, and Telugu in both directions. We used the Flores 200 devtest corpus and randomly selected 120 pairs of sentences. Three different raters were engaged to evaluate each pair of translations. The evaluated translation engines include Google Translate, IndicTrans2, GPT3.5, Llama-2-13b+FF+lora (Multi), and Llama-2-7b+lora(BI). The overall results for direct assessment scores (averaged on 120 sentences and 3 different ratings) are shown in Figure 6 for both translation directions. The overall ranking of different systems is similar to the automatic evaluation methods such as BLEU and CHRF scores. Our finetuned models on smaller parallel corpora for English-to-Indian-language machine translation perform better compared to GPT3.5. However, when we compare Indian Languages to English human evaluation, the performance is not the same. This is mainly attributed to the limited or near-zero vocabulary coverage in the LLM models. Furthermore, as already discussed, direct assessment indicates the superiority of encoder-decoder-based models for the translation task, such as IndicTrans2 and Google Translate.

Therefore, automatic and human evaluation suggest the need for large language models (LLMs) with sufficient representation for Indian languages. Future LLM development must address this require-



Figure 6: Human Evaluation comparison of Google Translate, Indictrans2, GPT-3.5 vs our Fine-Tuned LLM Translation models (LLaMA-2-13b+lora (BI), and LLaMA-2-13b+FF+lora (Multi)): English to 4 Indian languages on 120 sentences each with 3 ratings (averaged) and 4 Indian Languages to English on 120 sentences each with 3 ratings.

ment.

7 Conclusion

Our experiments and results have provided promising insights into the use of LLMs for translation tasks. We have found that LLMs have the potential to perform translations involving English and Indian languages without the need for an extensive collection of parallel data, which distinguishes them from traditional translation models. Furthermore, our findings indicate that the models based on LLaMA-2 outperform other models in the zero-shot and in-context example-based learning. In particular, the LLaMA-2-13b-based model demonstrates superior performance compared to its counterparts. To enhance the LLM’s understanding of English and Indian languages, we have introduced a two-stage fine-tuning process. This process begins with an initial full fine-tuning, followed by LoRa-based fine-tuning. Through this approach, we have significantly improved the LLM’s comprehension of content in both languages.

However, our experiments suggest that further work is required on LLMs to surpass the performance of traditional encoder-decoder-based translation models. This work could involve the development of LLMs specific to Indian languages, which would improve vocabulary and alphabet coverage, resulting in a better representation of Indian languages.

However, in the future, we plan to incorporate Indian-to-Indian language translation using LLM by exploring vocabulary expansion approaches. Furthermore, our objective is to develop a single LLM model capable of translating all Indian languages, as well as English, in both directions. By doing so, we aim to push the boundaries of language capabilities within LLMs and further advance the field.

8 Limitations

To conduct our experiments, we relied on high-performance GPUs, specifically the A100-40GB. However, we acknowledge that not everyone may have access to such powerful computing resources, making it challenging to reproduce our experiments and achieve identical results. To overcome this limitation, our objective is to provide open access to all outputs, including models and results, to facilitate further research and exploration. By making these resources openly available, we aim to promote collaboration and allow others to build on our work.

9 Ethics statement

To perform the experiments, we use publicly available data sets. Since we fine-tune the models on publicly available datasets, the models might not be prone to any ethical concerns. To encourage the reproducibility, we mention all the experimental details.

Acknowledgement

We express our gratitude to Pruthwik Mishra, Arafat Ahsan, and Palash Gupta for their input throughout the different phases of this project. This undertaking is funded by the Ministry of Electronics and Information Technology, Government of India, as evidenced by the Sanction Order: 11(1)/2022-HCC(TDIL)-Part(2)/A/B/C and the Administrative Approval: 11(1)/2022-HCC(TDIL)-Part(2).

References

AI4Bharat, Jay Gala, Pranjal A. Chitale, Raghavan AK, Sumanth Doddapaneni, Varun Gumma, Aswanth Kumar, Janki Nawale, Anupama

- Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M. Khapra, Raj Dabre, and Anoop Kunchukuttan. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *arXiv preprint arXiv:2305.16307*.
- Aksitov, Renat, Chung-Ching Chang, David Reitter, Siamak Shakeri, and Yunhsuan Sung. 2023. Characterizing attribution and fluency trade-offs for retrieval-augmented large language models. *arXiv preprint arXiv:2302.05578*.
- Aktar Husain, Jaavid, Raj Dabre, Aswanth Kumar, Ratish Puduppully, and Anoop Kunchukuttan. 2024. Romansetu: Efficiently unlocking multilingual capabilities of large language models via romanization. *arXiv e-prints*, pages arXiv–2401.
- Almazrouei, Ebtesam, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Coljocar, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malaric, et al. 2023. Falcon-40b: an open large language model with state-of-the-art performance. Technical report, Technical report, Technology Innovation Institute.
- Bang, Yejin, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.
- Bapna, Ankur and Orhan Firat. 2019. Exploring massively multilingual, massive neural machine translation. *Google AI Blog*, October, 11.
- Barrault, Loïc, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Yvette Graham, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, and Matteo Negri, editors. 2020. *Proceedings of the Fifth Conference on Machine Translation*, Online, November. Association for Computational Linguistics.
- Bellegarda, Jerome R. 2004. Statistical language model adaptation: review and perspectives. *Speech communication*, 42(1):93–108.
- Bengio, Yoshua, Réjean Ducharme, and Pascal Vincent. 2000. A neural probabilistic language model. *Advances in neural information processing systems*, 13.
- Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Chang, Yupeng, Xu Wang, Jindong Wang, Yuan Wu, Kaijie Zhu, Hao Chen, Linyi Yang, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2023. A survey on evaluation of large language models. *arXiv preprint arXiv:2307.03109*.
- Cui, Yiming, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*.
- Dabre, Raj, Chenhui Chu, and Anoop Kunchukuttan. 2020. A survey of multilingual neural machine translation. *ACM Computing Surveys (CSUR)*, 53(5):1–38.
- Eriguchi, Akiko, Shufang Xie, Tao Qin, and Hany Hassan Awadalla. 2022. Building multilingual machine translation systems that serve arbitrary xy translations. *arXiv preprint arXiv:2206.14982*.
- Federmann, Christian, Tom Kocmi, and Ying Xin. 2022. NTREX-128 – news test references for MT evaluation of 128 languages. In *Proceedings of the First Workshop on Scaling Up Multilingual Evaluation*, pages 21–24, Online, nov. Association for Computational Linguistics.
- Gala, Jay, Pranjal A Chitale, Raghavan AK, Sumanth Doddapaneni, Varun Gumma, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, et al. 2023. Indictrans2: Towards high-quality and accessible machine translation models

- for all 22 scheduled indian languages. *arXiv preprint arXiv:2305.16307*.
- Goyal, Naman, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.
- Guerreiro, Nuno M, Duarte Alves, Jonas Waldendorf, Barry Haddow, Alexandra Birch, Pierre Colombo, and André FT Martins. 2023. Hallucinations in large multilingual translation models. *arXiv preprint arXiv:2303.16104*.
- Holmström, Oskar and Ehsan Doostmohammadi. 2023. Making instruction finetuning accessible to non-English languages: A case study on Swedish models. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 634–642, Tórshavn, Faroe Islands, May. University of Tartu Library.
- Hu, Edward J, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Jiang, Albert Q, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Khan, Mohammed Safi Ur Rahman, Priyam Mehta, Ananth Sankar, Umashankar Kumaravelan, Sumanth Doddapaneni, Sparsh Jain, Anoop Kunchukuttan, Pratyush Kumar, Raj Dabre, Mitesh M Khapra, et al. 2024. Indicllmsuite: A blueprint for creating pre-training and fine-tuning datasets for indian languages. *arXiv preprint arXiv:2403.06350*.
- Kwon, Woosuk, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Lai, Viet Dac, Nghia Trung Ngo, Amir Poursan Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Huu Nguyen. 2023. Chatgpt beyond english: Towards a comprehensive evaluation of large language models in multilingual learning. *arXiv preprint arXiv:2304.05613*.
- Madaan, Pulkit and Fatiha Sadat. 2020. Multilingual neural machine translation involving Indian languages. In Jha, Girish Nath, Kalika Bali, Sobha L., S. S. Agrawal, and Atul Kr. Ojha, editors, *Proceedings of the WILDRE5–5th Workshop on Indian Language Data: Resources and Evaluation*, pages 29–32, Marseille, France, May. European Language Resources Association (ELRA).
- Mangrulkar, Sourab, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. 2022. Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>.
- Mehandru, Nikita, Samantha Robertson, and Nilofar Salehi. 2022. Reliable and safe use of machine translation in medical settings. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 2016–2025.
- Naveed, Humza, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Nick Barnes, and Ajmal Mian. 2023. A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Penedo, Guilherme, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023a. The

- refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.
- Penedo, Guilherme, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023b. The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.
- Popović, Maja. 2015. chrF: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels, October. Association for Computational Linguistics.
- Radford, Alec, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Rei, Ricardo, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Reinauer, Raphael, Patrick Simianer, Kaden Uhlig, Johannes EM Mosig, and Joern Wuebker. 2023. Neural machine translation models can learn to be few-shot learners. *arXiv preprint arXiv:2309.08590*.
- Salomon, Richard. 1995. On the origin of the early indian scripts. *Journal of the American Oriental Society*, 115(2):271–279.
- Scao, Teven Le, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Stanchev, Peter, Weiyue Wang, and Hermann Ney. 2020. Towards a better evaluation of metrics for machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 928–933.
- Steigerwald, Emma, Valeria Ramírez-Castañeda, Débora YC Brandt, Andrés Báldi, Julie Teresa Shapiro, Lynne Bowker, and Rebecca D Tarvin. 2022. Overcoming language barriers in academia: Machine translation tools and a vision for a multilingual future. *BioScience*, 72(10):988–998.
- Touvron, Hugo, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Touvron, Hugo, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Vaithilingam, Priyan, Tianyi Zhang, and Elena L Glassman. 2022. Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models. In *Chi conference on human factors in computing systems extended abstracts*, pages 1–7.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Wang, Rui, Xu Tan, Renqian Luo, Tao Qin, and Tie-Yan Liu. 2021. A survey on low-resource neural machine translation. *arXiv preprint arXiv:2107.04239*.
- Wu, Haoran, Wenxuan Wang, Yuxuan Wan, Wenxiang Jiao, and Michael Lyu. 2023a. Chatgpt or

- grammatically? evaluating chatgpt on grammatical error correction benchmark. *arXiv preprint arXiv:2303.13648*.
- Wu, Zhenyu, Yaoxiang Wang, Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Jingjing Xu, and Yu Qiao. 2023b. OpenICL: An open-source framework for in-context learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 489–498, Toronto, Canada, July. Association for Computational Linguistics.
- Xi, Zhiheng, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2023. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*.
- Xu, Haoran, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv preprint arXiv:2309.11674*.
- Xuanfan, Ni and Li Piji. 2023. A systematic evaluation of large language models for natural. In *Proceedings of the 22nd Chinese National Conference on Computational Linguistics (Volume 2: Frontier Forum)*, pages 40–56, Harbin, China, August. Chinese Information Processing Society of China.
- Zhang, Susan, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.
- Zhang, Biao, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study. *arXiv preprint arXiv:2301.07069*.
- Zhu, Wenhao, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian Huang. 2023. Multilingual machine translation with large language models: Empirical results and analysis. *arXiv preprint arXiv:2304.04675*.
- Zieliński, Jakub et al. 2021. Language as an interstate migration barrier—the interesting case of india. *Eastern European Journal of Transnational Relations*, 5(1):29–38.

A Encoder-Decoder BPCC (H) Model

We have trained a machine translation model using the human training data from BPCC in English to 22 Indian languages (AI4Bharat et al., 2023) on an A100 40GB GPU. The following are the details of the model configuration for training:

- Input: 32K merge operations-based subword tokens; Embedding size: 512, 4096 feedforward size; Layers: Encoder: 6, Decoder:6 and Attention heads: 8; Dropout: 0.30; Max word sequence length: 200; Steps: 200000; Batch Size: 8192 tokens; Initial learning rate: 2e-5; Optimizer: Adam; Label-smoothing: 0.1; 16-bit floating point precision; Early stop with no increase on training loss (10 epochs); Beam size: 15

B Large Language Models

Language modeling, a well-established task in the field of natural language processing, has attracted significant attention over the years (Bellegarda, 2004; Bengio et al., 2000). This task involves predicting the probability of the next token in a sequence of words. Transformers have emerged as the fundamental architecture underlying many existing large-language models (Vaswani et al., 2017).

Transformer-based autoregressive models, such as GPT (Brown et al., 2020; Radford et al., 2019) have played a crucial role in the advancement of natural language processing (NLP). GPT-3, with 175 billion parameters, is a standout in this category. It is similar in structure to GPT-2 and GPT-1 but benefits from a more extensive and varied dataset, making it exceptionally powerful in NLP. In addition, prompt-based ChatGPT (GPT-3.5 text-davinci-003 and GPT-3.5 turbo) has been performing exceptionally utilizing the reinforcement-based human feedback strategy. Although these models exhibit impressive performance on several NLP tasks, privacy and bias of the models have been a bottleneck. To mitigate such issues, LLaMA (Touvron et al., 2023a) is an open source foundation model trained on publicly available datasets. Similarly, Falcon-40B (Almazrouei et al., 2023) is another open-source LLM trained on a RefinedWeb corpus of 1500 billion tokens. Falcon even comes with 7 and 40 billion instruction versions trained on conversation data.

The recent adaptation of Large Language Models (LLMs) for instruction tuning has proven to be a

promising approach to improve the performance of various natural language processing tasks. Specifically, in languages like Chinese and Swedish, this shows the impressive zero-shot and generation abilities of the low-rank adaptation of LLaMA for non-English languages (Cui et al., 2023; Holmström and Doostmohammadi, 2023). The recent development of INDICLLMSUITE (Khan et al., 2024) is an initiative for large language models focusing on Indian languages. However, it is worth noting that the current focus of these instruction models is primarily on English. Therefore, there is an immediate need to explore ways to adapt these models to low-resource Indian languages (Aktar Husain et al., 2024).

B.1 Base/Raw Models

In this work, we use the following base LLM models to test the levels of language coverage and explore their potential for machine translation tasks involving English and Indian languages.

- **opt-6.7b**¹⁵ : The OPT-6.7b (Zhang et al., 2022) model has been extensively trained on the objective of causal language modeling (CLM) using English text. Although most of the training data are in English, a small portion of non-English data from CommonCrawl has also been included. This model utilizes 6.7 billion parameters, consisting of 32 layers and 32 attention heads, and employs an embedding size of 4096.
- **Bloom-7B**¹⁶ : BLOOM (Scao et al., 2022) was the first multilingual large language model with a causal language modeling objective and supports 46 languages and 13 programming languages. Its overall training data contains 1.1% of Indian languages. We opted for Bloom model with 7,069,016,064 parameters with 30 layers, 32 attention heads, 4096 embedding dimensional where the maximum token length is 2048.
- **LLaMA-7B**¹⁷: LLaMA is a collection of foundation language models that range from 7B to 65B parameters. These models are

¹⁵<https://huggingface.co/facebook/opt-6.7b>

¹⁶<https://huggingface.co/bigscience/bloom-7b1>

¹⁷<https://huggingface.co/decapoda-research/llama-7b-hf>

multilingual models and are trained on trillions of tokens. Data include CCNet, C4, GitHub, Wikipedia, Books, ArXiv, and Stack Exchange. In our experiments, we evaluated the LLaMA model with 7B parameters where 4096 is the embedding dimension and 32 layers and 32 attention heads.

which reduces the memory requirement during decoding. It has 4096 embedding dimensions, 32 layers, and 32 attention heads with context length of 8192 context length.

- **MPT-7B**¹⁸ : Similarly to the above models, the MPT-7B model is trained on a large number of 1T data tokens in the causal language modeling objective.
- **Falcon**¹⁹ : Falcon (Penedo et al., 2023b) is another large language model trained on causal language modeling (CLM). Here, we utilized Falcon-7B model which is a 7B parameters and trained on 1.5 trillion tokens of Refined-Web (a novel massive web data set based on CommonCrawl) enhanced with curated corpora. The model has multilingual capabilities, but Indian languages are not explicitly present. We used Falcon-7B for our experiments.
- **LLaMA-2-7B**²⁰ and **LLaMA-2-13B**²¹ : LLaMA 2 based models (Touvron et al., 2023b) are also trained on the causal language modeling (CLM) objective and pretrained on 2 trillion tokens of data from publicly available sources of up to September 2022. These models are available in different range parameters from 7 billion to 70 billion. These models have 4k subwords as the context length. In our experiments, we have experimented with 7B and 13B LLaMA-2 models. LLaMA-2-7B network has 32 layers and 32 attention heads, while the LLaMA-2-13B network has 40 layers and 40 attention heads.
- **Mistral-7B**²² : Mistral-7B Large Language Model (LLM) (Jiang et al., 2023) is a pre-trained with causal language modeling (CLM) objective with 7 billion parameters. It uses Sliding-Window Attention (SWA) to handle longer sequences at a lower cost and grouped query attention (GQA) for faster inference,

¹⁸<https://huggingface.co/mosaicml/mpt-7b>

¹⁹<https://huggingface.co/tiiuae/falcon-7b>

²⁰<https://huggingface.co/meta-llama/Llama-2-7b-hf>

²¹<https://huggingface.co/meta-llama/Llama-2-13b-hf>

²²<https://huggingface.co/mistralai/Mistral-7B-v0.1>

C MT systems outputs

We have added some examples of our best-performing models for zero-shot, ICL, fine-tuning and 2-stage fine-tuning strategies in Figures 7, 8, 9, 10, 11, 12.

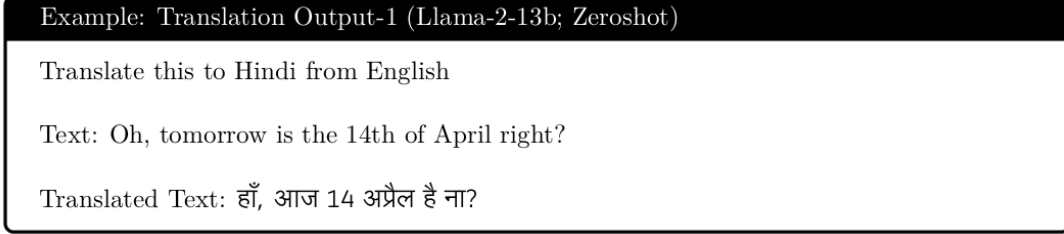


Figure 7: Translation example for Llama-2-13b model with zero-shot setting.

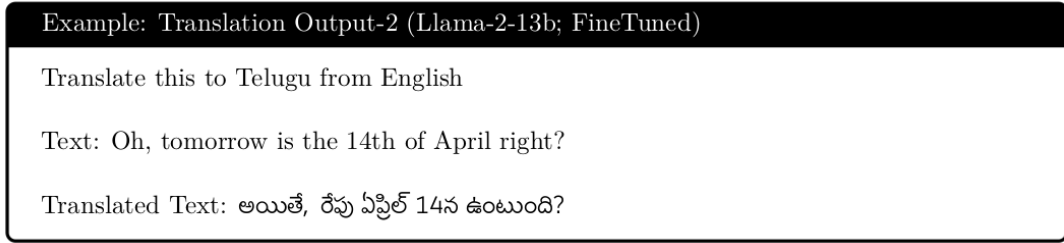


Figure 8: Translation example for Llama-2-13b finetuned model.

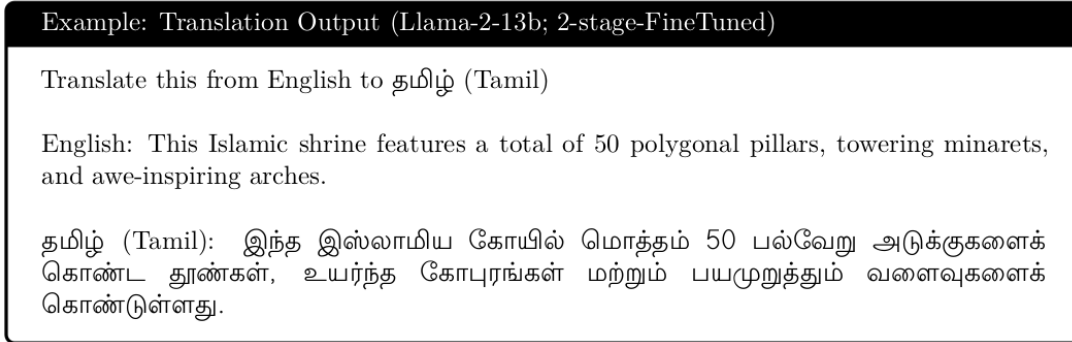


Figure 9: Translation example for Llama-2-13b 2-stage-finetuned model.

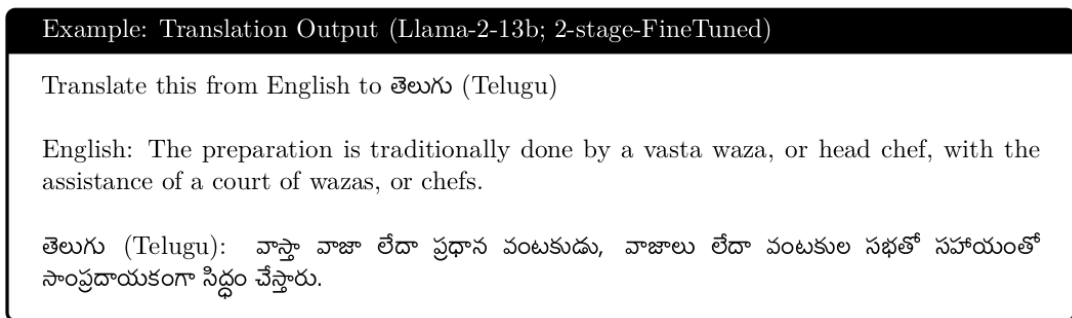


Figure 10: Translation example for Llama-2-13b 2-stage-finetuned model.

Example: Translation Output with ICL (Llama-2-7b)

If the English to Hindi translation for 'Charbaugh is the railway station while Alambaugh and Kaisarbaugh are major bus terminus here , from where local travel resources like taxis - scooter rickshaw etc. can be used for tour .' is 'चारबाग रेलवे स्टेशन है , जबकि आलमबाग और कैसरबाग यहाँ के प्रमुख बस अड्डे हैं , जहाँ से सैर के लिए स्थानीय यातायात के साधनों मसलन टैक्सियों – स्कूटर रिक्शा आदि का इस्तेमाल किया जा सकता है ।', following that, translate this to Hindi from English.

Text: Are all schools, colleges and offices closed tomorrow?

Translated Text: क्या हर स्कूल , कॉलेज और कार्यालय के लिए कल बंद होंगे ?

Figure 11: Translation example for Llama-2-7b In-context learning strategy.

Example: Translation output with hallucinations (Llama-2-13b-Finetuned)

Translate this to Hindi from English

Text: That's a good news!

Translated Text: यह अच्छा खबर है!

Text: That's a bad news!

Translated Text:यह खराब खबर है!

Text: That's a great news!

Translated Text:यह बहुत अच्छा खबर है!

Text: That's a terrible news!

Translated Text:यह बहुत खराब खबर है!

Figure 12: Hallucinated translation example for Llama-2-13b 2-stage-fine-tuned model.

D Detailed experimental results

We have listed the detailed experimental results in for all the 22 languages, with various combinations of data sets and models in Tables 5, 6, 7, 8.

Improving NMT from a Low-Resource Source Language: A Use Case from Catalan to Chinese via Spanish

Yongjian Chen¹, Antonio Toral¹, Zhijian Li², Mireia Farrús^{3,4}

Center for Language and Cognition, University of Groningen, Netherlands¹

School of Foreign Languages, Guangzhou City University of Technology, China²

Centre de Llenguatge i Computació, Universitat de Barcelona, Spain³

Institut de Recerca en Sistemes Complexos, Universitat de Barcelona, Spain⁴

{yongjian.chen, a.toral.ruiz}@rug.nl¹

lizhijian@geu.edu.cn²

mfarrus@ub.edu^{3,4}

Abstract

The effectiveness of neural machine translation is markedly constrained in low-resource scenarios, where the scarcity of parallel data hampers the development of robust models. This paper focuses on the scenario where the source language is low-resource and there exists a related high-resource language, for which we introduce a novel approach that combines pivot translation and multilingual training. As a use case we tackle the automatic translation from Catalan to Chinese, using Spanish as an additional language. Our evaluation, conducted on the FLORES-200 benchmark, compares our new approach against a vanilla baseline alongside other models representing various low-resource techniques in the Catalan-to-Chinese context. Experimental results highlight the efficacy of our proposed method, which outperforms existing models, notably demonstrating significant improvements both in translation quality and in lexical diversity.

1 Introduction

The development of neural machine translation (NMT) has considerably benefited translation between language pairs abundant in parallel data, enhancing translation accuracy and fluency across diverse linguistic landscapes (Hassan Awadalla et al., 2018; Popel et al., 2020). However, its effect is challenged by the fact that building an effective NMT system requires a large amount of

parallel data. This challenge is particularly pronounced in the case of low-resource languages, that is, language pairs with limited parallel language resources, remaining a significant hurdle in achieving universal communication.

An exemplary case highlighting such hurdle involves the translation dynamics between Catalan and Chinese, (CA–ZH) two languages characterized by limited parallel corpora. The year 2022 marked a significant increase in Chinese investments in Catalonia¹, and the Chinese population emerged as the fourth largest foreign community in Catalonia². These together highlight the growing economic and social interactions between these regions and thus the pressing need for effective communication tools between Catalan and Chinese speakers. Despite the potential benefits, the development of a robust NMT system for the CA–ZH language pair faces notable challenges, primarily due to the scarcity of direct parallel data.

Addressing this gap, previous works have sought to navigate the low-resource landscape of the CA–ZH language pair. The research by Costa-Jussà et al. (2019) was the first work to specifically focus on addressing the low-resource CA–ZH language pair, where they broke new ground by generating non-human-written parallel sentences, i.e. pseudo-parallel corpus via pivot translation and then used them to train ZH→CA NMT models. Another work (Zhou, 2022) concerned building CA–ZH parallel data, where CA–ZH bitexts was first mined from Wikipedia with the help of the open-source *LASER* toolkit³ and then passed to san-

¹Data taken from <https://catalonia.com>.

²Data taken from <https://www.idescat.cat/novetats/?id=4489&lang=en>.

³<https://github.com/facebookresearch/LASER>

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

ity check according to Kreutzer’s (2022) methodology. Therefore, unlike the training datasets created by Costa-Jussà et al. (2019), Zhou’s (2022) parallel corpus consists of human-selected bitexts. Subsequently, Liu (2022) used Zhou’s (2022) dataset to fine-tune a massively pre-trained multilingual NMT model, i.e. M2M-100 with parameters of 418M by Fan et al. (2020) for CA↔ZH, presenting better translation performance in both directions as compared to the original M2M-100.

Furthermore, other work (Schwenk et al., 2019; El-Kishky et al., 2020; Schwenk et al., 2020), related to building parallel data for many language pairs, included CA–ZH. The multilingual bitexts therein were massively mined from web-based resources and subsequently utilized to train multilingual NMT models like M2M-100, which is also workable with the CA↔ZH translation.

The previous studies on the CA–ZH language pair contributed to enhancing the automated translation between these two languages. However, each focused primarily on employing a singular, specific low-resource NMT technique, e.g. pivot translation in Costa-Jussà et al.(2019), fine-tuning in Liu (2022), multilingual training in Fan et al. (2020), etc. Unlike these approaches, our work aims to propose a novel integration, pivot translation-aided multilingual training (PTAMT), and compare it against existing methods (multilingual training, fine-tuning, pivot translation, fine-tuning coupled with pivot translation). We focus on the CA→ZH translation direction, as a use case in which the source language is low-resource and there is a related higher-resourced language, Spanish (ES). The technique we introduce, PTAMT, uses additional data from ES both as pivot and as multilingual training.

The contributions of our work can be summarized as follows:

1. Our work introduces a novel approach that effectively leverages pseudo-parallel and authentic data to enhance translation quality and mitigate the effects of source-side *machine translationese*, setting a new standard for NMT from low-resource languages.
2. Our work, to the best of our knowledge, is the first one to provide systematic empirical evidence highlighting the effectiveness of different low-resource NMT techniques for the CA–ZH language pair.

3. Our work underscores the important role of a modest amount of authentic parallel data in the target language pair(s) in the training and fine-tuning processes.

2 Related Work

2.1 Low-resource Techniques

Multilingual training refers to training for different language pairs in a single NMT model (Wang et al., 2021) via various methods of sharing parameters, e.g. full parameters sharing (Ha et al., 2016; Johnson et al., 2017; Tan et al., 2019), attention mechanism sharing (Firat et al., 2016; Lu et al., 2018), etc. Through multilingual training, low-resource language pairs can be trained together with high-resource language pairs, and thus desired low-resource languages can benefit from high-resource auxiliary languages when the model learns linguistic knowledge, contextual information, and commonalities, etc. from different languages. Furthermore, if auxiliary languages are related to low-resource languages of interest, they can effectively benefit translation quality in a low-resource scenario (Gu et al., 2018; Neubig and Hu, 2018).

Fine-tuning is performed when a parent NMT model is first trained on high-resource language pairs, and the trained model is used to initialize a child model’s parameters, which are subsequently fine-tuned on a low-resource language pair (Zoph et al., 2016). In this way, whereas knowledge learnt from high-resource auxiliary languages can be transferred to low resource languages, the pre-trained NMT model can also be forced to primarily focus on the desired low-resource language pair only. By contrast, since a model has constrained capacity, multilingual training can potentially favor high-resource language pairs due to imbalanced data ratio (Arivazhagan et al., 2019; Wang et al., 2020). Fine-tuning can be combined with multilingual training if a model is first trained on multiple high-resource languages as well as the desired low-resource language pair and then is fine-tuned on the latter only, which has been proved as an effective way to improve low-resource translation (Thillainathan et al., 2021; Adelani et al., 2022).

Pivot translation is applicable for a low-resource translation condition if an auxiliary language has parallel data with both languages of the

low-resource language pair, and this auxiliary language is called a pivot language (Costa-Jussà et al., 2019).

There are mainly two approaches in pivot translation. The first one is the cascade approach, aiming to train two separate translation systems from source to pivot and from pivot to target, and then combine them together for source→target translation, which is common in early statistical machine translation (Cohn and Lapata, 2007; Wu and Wang, 2007; El Kholly et al., 2013).

Another approach is more widely used in state-of-the-art NMT, which is used to synthesize pseudo-parallel data for a low-resource language pair, with data either from the source side synthesized through pivot→source translation (Zheng et al., 2017) or from the target side synthesized through pivot→target translation (Chen et al., 2017). In this case, to ensure the effectiveness of synthetic data via pivot translation, it is important to obtain qualified pivot→source translation or pivot→target translation. For instance, Costa-Jussà et al. (2019) compared two pseudo-parallel CA–ZH parallel corpora. One was built by translating the Spanish sentences from the ES–ZH parallel corpus *United Nations Parallel Corpus v1.0* (Ziems et al., 2016) into Catalan, whereas the other was created by translating the Spanish sentences from the ES–CA parallel corpus *El Periódico* (Costa-jussà et al., 2014) into Chinese. They used them to train two separate NMT models with same neural network architecture for ZH→CA translation, and discovered that the NMT model trained on the former yielded a higher BLEU score, as the ES→CA translation was of higher quality than the ES→ZH translation.

In contrast to this direct synthesis approach, studies by Lakew et al. (2018) and Currey and Heafield (2019) leveraged pivot resources differently. Rather than solely relying on pivot→source or pivot→target translations to generate pseudo-parallel data, these studies initiated their process with multilingual NMT model training using both source–pivot and target–pivot parallel data. Afterwards, Lakew et al. (2018) used their multilingual NMT model to back-translate source and target language data into the corresponding target and source languages, respectively. This generated pseudo-parallel source–target data was then used alongside the original parallel data to iteratively re-train the multilingual NMT model. Differ-

ently, Currey and Heafield (2019) leveraged their multilingual model to back-translate monolingual data from the pivot language into both the source and target languages, thereby obtaining pseudo-parallel source–target data used to further train or fine-tune the model to enhance the direct translation capabilities between the source and target languages.

2.2 Machine Translationese

Machine translationese refers to the artificially impoverished language in MT outputs, marked by reduced fluency, lexical diversity, and distinct syntactic structures compared to original or human-translated texts (Vanmassenhove et al., 2021; Chae and Nenkova, 2009; Ilisei et al., 2010). Such characteristics can make synthetic machine translated data ill-suited for capturing the nuances of human language, potentially leading to deviations in real-world language usage (Dutta Chowdhury et al., 2022). When synthetic data is utilised as training data (as can be the case in pivot translation, see Section 2.1), models may inadvertently learn from the machine translationese present in the synthetic data, leading to the generation of translations or language constructs that are inconsistent with the target language.

Efforts to mitigate translationese have included techniques such as data tagging, where training datasets are annotated to distinguish between original and translated texts. This tagging helps models recognize and avoid translationese during training, as in Riley et al. (2020) and Freitag et al. (2022b). Another approach involves transforming machine-translated texts into more original-like content using style transfer or by re-generating text from abstracted representations like AMR (Jalota et al., 2023; Wein and Schneider, 2024).

These studies mainly focus on improving the quality of machine-translated output by reducing translationese. However, less attention has been given to the effects of source-side artefacts in synthetic data on NMT training. We contemplate this case in this work, comparing different models that deal with synthetic source-side training data in terms of machine translationese.

3 Proposed Method

Our proposed method, PTAMT, couples pivot translation with multilingual training to leverage the advantages of both techniques. Distinct from

previous work (Lakew et al., 2018; Currey and Heafield, 2019), which uses a less related pivot language to initially train a multilingual NMT system for back-translating and synthesizing pseudo-parallel data, PTAMT employs a pivot language (ES) that is linguistically closer to the source language. This choice is informed by the synthetic pivot translation approach demonstrated by Costa-Jussà et al. (2019), which is favored due to the greater linguistic affinity between the source (CA) and pivot (ES) languages as compared to the pivot (ES) and target (ZH) languages (Rapp, 2021).

In our implementation, we used an existing ES–ZH corpus to synthesize pseudo-parallel CA–ZH data by translating ES sentences into CA. This strategy aligns the synthetic side with the source language (CA) and uses authentic data for the target language (ZH), enhancing the model’s ability to produce natural output.

Nevertheless, as pivot-translated sentences are inherently machine-translated texts, they are prone to containing machine translationese. To address this, PTAMT strategically leverages the ES–ZH bitexts in a multilingual training setup to facilitate effective pivot-based knowledge transfer. This approach helps to mitigate the potential impact of noise introduced by the synthetic CA input. PTAMT incorporates both source languages (CA and ES) in the encoder, and applies encoder parameter sharing throughout training, which is applicable for both from-scratch training and fine-tuning. The same set of weights and biases is shareable in a single encoder for inputs from both source languages. Multilingual training empowers the model to capture and integrate contextual information from both source languages. Given that the ES data is authentic (human-produced), we hypothesize that PTAMT will aid in reducing the influence of noise from the pivot-translated CA training data on the target ZH output.

4 Experimental Setup

4.1 Data Description

Original Data We used the aforementioned CA–ZH parallel corpus, *CA–ZH Wikipedia* (Zhou, 2022), as the foundation, since it contains human-selected parallel sentences with quality control. We refer to this dataset as *CA–ZH-WIKI*.

Pivot-translated Data We made use of ES–ZH bitexts from the public release *United Nations*

Parallel Corpus v1.0 (Ziemski et al., 2016), to which we refer as *ES-ZH-UN*. Instead of training an ES→CA automatic translation system from scratch as in (Costa-Jussà et al., 2019), we directly used the open-source NMT model provided by Softcatalà⁴ to translate ES sentences from *ES-ZH-UN* into CA and then obtain the pseudo-parallel (synthetic source) CA–ZH United Nations parallel dataset, *CA-ZH-PVT*. Additionally, by translating CA sentences from *CA-ZH-WIKI* into ES through a CA→ES NMT model⁵, we generated a pseudo ES–ZH Wikipedia parallel dataset, *ES-ZH-PVT*, for later data augmentation.

Mixed Data As for the CA–ZH language pair, we concatenated the pivot-translated CA–ZH dataset *CA-ZH-PVT* with *CA-ZH-WIKI*, resulting in the mixed parallel dataset *CA-ZH-MIX*. As regards the ES–ZH language pair, we concatenated the pivot-translated ES–ZH Wikipedia dataset *ES-ZH-PVT* with *ES-ZH-UN* to obtain a mixed parallel dataset *ES-ZH-MIX* for the ES–ZH language pair.

4.2 Models

We implemented one vanilla baseline model, four models based on a pre-trained multilingual model (M2M-100-418), two Transformer-based models trained from scratch, and four PTAMT models to compare different low-resource NMT techniques for CA→ZH translation.

Vanilla Baseline Our vanilla baseline was a Transformer-based model trained on the original parallel corpus *CA-ZH-WIKI*. Due to the small size of this training dataset, rather than using the default Transformer-base configuration (Vaswani et al., 2017), we adopted the architecture setting optimized on 40k training sentence pairs (Araabi and Monz, 2020), which consists of 2 attention heads, 5 encoder and decoder layers, and a 512 embedding dimension.

M2M-100-418M Models As a second baseline, we selected the pretrained model M2M-100 (Fan et al., 2020), which is representative of a model that has taken advantage of multilingual training. M2M-100 is a state-of-the-art massively multilingual translation model, which supports translation between Catalan and Chinese. We opted for the

⁴<https://github.com/Softcatala/nmt-models>

⁵<https://github.com/Softcatala/nmt-models>

Language pair	Corpus	# of sentence pairs	
		Training	Validation
CA-ZH	CA-ZH-WIKI	58,328	10,293
	CA-ZH-PVT	17,575,795	2,638
	CA-ZH-MIX	17,634,123	12,931
ES-ZH	ES-ZH-UN	17,575,795	2,638
	ES-ZH-PVT	58,328	10,293
	ES-ZH-MIX	17,634,123	12,931

Table 1: Distribution of the datasets in the experiments.

one with the least size of parameters (418M) taking into account available computational resources as well as comparability across the different models in our experiments. M2M-100-418M is a Transformer-based model that contains 12 encoder and decoder layers with a 1024 embedding dimension.

Large-scale multilingual pre-trained NMT models can be further leveraged to improve low-resource machine translation by fine-tuning them on low-resource language pairs. Therefore, we examined three fine-tuned M2M-100-418M models for the CA→ZH translation. The first one was obtained from the aforementioned work by Liu (2022), accessible on a Hugging Face repository⁶, which was solely fine-tuned on the *CA-ZH-WIKI* training dataset. We fine-tuned the second one on the pseudo-parallel dataset *CA-ZH-PVT* and the third one on the mixed parallel dataset *CA-ZH-MIX*. These last two models represent those that leverage pivot translation (either without or with original parallel data) paired with fine-tuning.

From-scratch Trained Models We additionally trained two models from scratch using a Transformer architecture, with 6 encoder and decoder layers and a 512 embedding dimension, respectively on the pseudo-parallel dataset *CA-ZH-PVT* and the mixed parallel dataset *CA-ZH-MIX*. These two models represent those that leverage pivot translation (either without or with original parallel data) under from-scratch training conditions.

PTAMT-enhanced Models We implemented PTAMT to enable simultaneous benefits from *CA-ZH-MIX* and *ES-ZH-MIX* in both from-scratch training and fine-tuning scenarios. Under the from-scratch training condition, we trained a sin-

gle NMT model that has the same network architecture as the previous from-scratch trained models. Whereas the language pair of interest is still CA→ZH, this model supports both CA→ZH and ES→ZH translation, effectively operating as a many-to-one NMT system. The encoder parameters are shared between CA and ES without increasing the model size, where a special token was added to the source side to specify the input language. Likewise, in the fine-tuning condition, both language pairs were included, and thus M2M-100-418M was fine-tuned on both *CA-ZH-MIX* and *ES-ZH-MIX*.

During the training or fine-tuning phase, ES was engaged as an auxiliary language. Despite potential noise introduced by pivot-translated CA sentences, the model could still learn relevant linguistic properties and characteristics related to CA from their ES equivalents, and thereby enhancing the CA→ZH translation. These two models represent PTAMT in from-scratch training and fine-tuning scenarios, respectively. Furthermore, we applied a second-step fine-tuning to both models on *CA-ZH-WIKI*.

4.3 Preprocessing

As for the parallel datasets used in our experiments, we only worked on sentence-level translation and so we removed lengthy sentence pairs by restricting them to maximum length of 100 words, then split them into training set and validation set (see Table 1), and went through different preprocessing pipelines depending on the models to be trained, as detailed next.

M2M-100-418M models For this model and its fine-tuned variants, including two with PTAMT, we employed the pre-trained SentencePiece tok-

⁶https://huggingface.co/projecte-aina/m2m100_418M_ft_ca_zh

enizer designed for M2M-100⁷. This tokenizer was used to tokenize all the parallel sentences.

Other models For the remaining models, we applied pre-tokenization using separate segmenters tailored for each language, following the approach outlined in Costa-Jussà (2019):

- **Chinese:** Since word boundaries in Chinese are not discernible through whitespace, we utilized the Jieba segmenter⁸ to segment Chinese sentences into words.
- **Catalan and Spanish:** We relied on the spaCy tokenizer, specifically the models *ca_core_news_sm*⁹ and *es_core_news_sm*¹⁰, respectively. These models were used to identify word boundaries and split contractions (e.g., *l'original* into *l' + original*).

Following the pre-tokenization step, we trained SentencePiece BPE models using training sentences from the respective datasets and then proceeded with tokenization.

- **Vanilla Baseline:** Following Araabi and Monz (2020), we trained a tokenizer with 12k BPE merge operations for each language.
- **From-scratch Models:** We sampled 5M sentences from the corresponding training set for each language. We then trained a tokenizer with a character coverage of 1.0 for CA and another one for ZH with character coverage of 0.9995. To determine the optimal vocabulary size for training our tokenizers, we initially used the widely-adopted size of 32k. Subsequently, we conducted experiments by both increasing and decreasing the vocabulary size. In evaluating the performance of tokenizers with different sizes, we assessed the occurrences of the *unk* token in the tokenized data. This resulted in a vocabulary of 35K subwords.
- **PTAMT:** For the proposed PTAMT-enhanced model under the from-scratch condition, we sampled 5M sentences from the ES corpus and concatenated it with the CA samples.

⁷https://dl.fbaipublicfiles.com/m2m_100/spm.128k.model

⁸<https://github.com/fxsjy/jieba>

⁹<https://spacy.io/models/ca>

¹⁰<https://spacy.io/models/es>

This combined set was used to train a joint tokenizer for CA and ES. After testing different vocabulary sizes, we finally created a joint vocabulary of 64K. We retained the Chinese tokenizer used in the previous from-scratch trained model for tokenization.

4.4 Training

For maximum comparability across the various models in this study, we conducted all experiments using a single NVIDIA RTX A5000 GPU card. We trained or fine-tuned all models with the Adam Optimizer and label smoothing cross-entropy loss. The configuration of hyper-parameters for all the NMT models is provided in Table 4 (Appendix A), except that hyper-parameters for training the vanilla baseline followed the ones in Araabi and Monz (2020) (see Table 5 in Appendix A). Additionally, checkpoints were evaluated at an intervals 5k training or fine-tuning steps on the validation set. Throughout this process, we continuously monitored the models' performance by assessing both training and validation losses. To ensure a balance between achieving convergence and avoiding overfitting, we implemented early stop if there was no improvement in the validation loss over 0.02 across three consecutive validation intervals. The epochs are listed in Table 6 in Appendix B.

4.5 Evaluation

4.5.1 Evaluation Benchmark

We benchmarked the models in this work on *FLORES-200* (Team et al., 2022). We used 1012 sentence pairs from its *devtest* set to evaluate the translation quality in the CA→ZH direction in all experiments, where we performed beam search decoding with a beam size of 5.

4.5.2 Evaluation Metrics

We incorporated three distinct sets of automatic evaluation metrics, with the first two aiming to evaluate translation quality and the last one aiming to assess lexical diversity.

SentencePiece BLEU We adopted the SentencePiece BLEU (spBLEU) (Goyal et al., 2021) as one of our quality evaluation metrics, since spBLEU correlated with human ratings slightly better than BLEU (Freitag et al., 2022a). We first detokenized the output from all the NMT models, then imple-

System	Methods	spBLEU	COMET
Vanilla Baseline	-	8.2	0.525
M2M-100-418M Models	multilingual training (pre-trained baseline)	22.0	0.774
	fine-tuning	22.4	0.797
	fine-tuning & pivot (without original data)	22.7	0.779
	fine-tuning & pivot (with original data)	24.6	0.808
	PTAMT	25.2	0.810
	PTAMT & 2nd fine-tuning	26.7	0.828
From-scratch Trained Models	pivot (without original data)	19.8	0.738
	pivot (with original data)	21.1	0.763
	PTAMT	23.1	0.783
	PTAMT & 2nd fine-tuning	24.3	0.786

Table 2: Translation quality automatic scores for the baseline, pre-trained models and from-scratch models. The best score per section and metric is shown in bold.

mented the pre-trained SentencePiece tokenizer¹¹ specific for FLORES-200 to tokenize the MT output and the reference translation, and finally computed spBLEU for each model.

Crosslingual Optimized Metric for Evaluation of Translation COMET (Rei et al., 2020) leverages cross-lingual neural language modelling and is trained to predict human judgement scores for machine-translated texts. COMET caters for a great variety of languages, and takes into account semantic similarities not only between the MT output and the reference translation but also the corresponding source text (Rei et al., 2020). We used the default COMET model¹², feeding it a triplet with detokenized source, MT output, and reference translation.

Measures of Lexical Diversity As discussed in Section 2.2, machine-translated texts exhibit differences in lexical diversity compared to original texts. Therefore, we evaluated lexical diversity in both reference translations and outputs from the NMT models in our experiments to compare the prevalence of machine translationese. Following the approach outlined by Vanmassenhove et

al. (2021), lexical diversity was examined using various measures, including lexical frequency profile (LFP), type/token ratio (TTR), Yule’s I and the measure of textual lexical diversity (MLTD).

In Vanmassenhove et al. (2021), LFP is used to quantify the richness of a translation by dividing the words of a text into three bands: (i) the percentage of words among the 1000 most common words in that language, (ii) the percentage of words among the next 1000 most common words, and (iii) all other words. These word frequency lists are generated from the training set. TTR assesses a text’s repetitiveness by comparing the ratio of unique words (types) to the total number of words (tokens) in the text. MLTD represents the mean length of a text where a given TTR value is maintained. Yule’s I, the inverse of Yule’s K, measures the constancy of text and the repetitiveness of vocabulary.

Prior to computing the lexical diversity scores for each metric, we tokenized the Chinese references and MT outputs following the same Chinese pre-tokenization steps outlined in Section 4.3. Besides, we utilized the pre-tokenized mixed Chinese training sentences from *CA-ZH-MIX* to obtain the Chinese word frequency list for LFP.

¹¹<https://github.com/facebookresearch/flores/tree/main/flores200>

¹²<https://huggingface.co/Unbabel/wmt22-comet-da>

System	Methods	B1	B2	B3	TTR	Yule’s I	MLTD
Reference	-	0.487	0.090	0.423	0.320	14.679	218.133
Vanilla Baseline	-	0.593	0.074	0.333	0.228	3.873	50.784
M2M-100-418M Models	multilingual training (strong baseline)	0.534	0.089	0.377	0.248	6.091	115.470
	fine-tuning	0.519	0.092	0.389	0.282	8.916	132.883
	fine-tuning & pivot (without original data)	0.517	0.093	0.391	0.279	10.046	121.555
	fine-tuning & pivot (with original data)	0.514	0.094	0.393	0.283	10.257	129.594
	PTAMT	0.517	0.093	0.390	0.288	10.770	157.991
	PTAMT & 2nd fine-tuning	0.515	0.088	0.396	0.295	10.011	167.163
From-scratch Trained Models	pivot (without original data)	0.584	0.095	0.321	0.246	6.762	137.691
	pivot (with original data)	0.579	0.091	0.331	0.256	6.927	127.027
	PTAMT	0.556	0.093	0.351	0.272	8.515	155.443
	PTAMT & 2nd fine-tuning	0.547	0.090	0.363	0.274	7.728	149.375

Table 3: LFP scores with 3 bands (B1: 0-1000, B2: 1001-2000, B3: 2001-end), TTR, Yule’s I and MLTD scores for the reference and the output of the NMT models in CA→ZH translation. Lower B1 values, indicating fewer matched tokens in frequent cases, along with higher values in B3, TTR, Yule’s I, and MLTD, collectively indicate greater lexical richness.

5 Results and Discussion

5.1 Results

Table 2 displays the quality outcomes, while Table 3 shows the lexical diversity outcomes for all the NMT models involved in the CA→ZH translation on the FLORES-200 dataset.

Translation Quality The results indicate a notable advancement in translation quality when examining the from-scratch trained models. Particularly, transitioning from the vanilla baseline to a pivot strategy yields a significant increase in performance metrics, with spBLEU surging by 11.6 points from 8.2 to 19.8, and the COMET score enhancing by approximately 0.213 from 0.525 to 0.738. This trend of improvement extends when integrating the pivot-translated dataset with the original, which further elevates the spBLEU score by 1.3 to 21.1. This enhancement is surpassed by the PTAMT model, marking a spBLEU increase of 2.0 from 21.1 to 23.1. Interestingly, fine-tuning the PTAMT model on the small amount of original dataset led to a further spBLEU boost by 1.2.

In comparison, the M2M-100-418M models begin with a strong foundation, exhibiting a high initial spBLEU score of 22 and a COMET score of 0.774. A slight improvement in spBLEU is noted

after fine-tuning on the original training set, increasing modestly to 22.4. The incremental advancement persists when pairing fine-tuning with pivot translation, further elevating the spBLEU to 22.7 when excluding the original parallel data and to 24.6 when combined with the original parallel dataset. Applying PTAMT in the fine-tuning condition boosts spBLEU further to 25.2, with a second-step fine-tuning on the original dataset resulting in a peak spBLEU score of 26.7, accompanied by the highest COMET score of 0.828.

The M2M-100-418M models generally outperform from-scratch models in terms of translation quality. However, the PTAMT-enhanced model in the from-scratch training scenario, whether with second-step fine-tuning or not, still surpasses the M2M-100-418M models reliant solely on multilingual training, fine-tuning, and fine-tuning combined with pivot translation (without original parallel data) in terms of spBLEU scores.

Lexical Diversity Compared to all the NMT models, the reference translation exhibits a lower B1 score and a higher B3 score. This reveals that the 1000 most frequent words represent a smaller proportion of the human-translated sentences, while less frequent words constitute a

larger portion of the original data compared to the outputs of different NMT systems, indicating a preference for less frequent words and a richer vocabulary. This is further evidenced by its superior TTR, Yule’s I, and MLTD scores.

However, the results also reveal that incorporating low-resource training approaches into NMT models consistently leads to performance improvements over the vanilla baseline. Specifically, our proposed method, PTAMT, stands out in both from-scratch trained models and atop the M2M-100-418M pre-trained model, by achieving the lowest B1 score, the highest B3 score, and the highest scores of TTR, Yule’s I, and MLTD. This suggests that the PTAMT-enhanced models excel in generating linguistically rich and varied outputs across both from-scratch training and fine-tuning scenarios. Furthermore, it was observed that the M2M-100-418M models demonstrate a superior ability to use a wider vocabulary compared to the from-scratch trained models.

5.2 Discussion

Notable improvements in translation quality and lexical diversity have been observed following the implementation of low-resource NMT techniques, underscoring the pivotal role of innovative training strategies in surpassing the limitations traditionally associated with NMT models in low-resource contexts.

While different approaches have exhibited different degrees of enhancement, the overall superiority of the M2M-100-418M models can be attributed to the extensive multilingual pre-training of the initial M2M-100-418M model, which is equipped with a broad variety of linguistic knowledge, enabling itself to benefit substantially from subsequent low-resource training strategies. Among the low-resource methods examined, PTAMT has set a new standard for generating translations, allowing the M2M-100-418M model to capitalize on both the CA–ZH and ES–ZH training datasets, achieving superior translation quality and lexical diversity compared to the other M2M-100-418M models examined. Interestingly, despite the inherent advantages of the M2M-100-418M’s large-scale multilingual training base, the from-scratch trained models leveraging the PTAMT method exhibit unique capacity to optimize translation quality beyond the capabilities of the M2M-100-418M models that rely

solely on the approaches of multilingual training, fine-tuning, and fine-tuning combined with pivot translation (without original parallel data).

Furthermore, PTAMT is particularly effective in reducing the impact of source-side machine-translationese introduced by the pivot-translated data (i.e. source-side machine-translated Catalan from Spanish) on the target output. PTAMT does not only include as training data a large amount of pseudo-parallel data for the desired source-target language pair (CA–ZH), but also integrates authentic linguistic input from the pivot–target language pair (ES–ZH). This approach enables the models to not only be exposed to a wider range of lexical items and usage contexts but also effectively discern and replicate the subtleties of natural language usage. Empirical evidence from our results of the lexical diversity metrics corroborates PTAMT’s positive impact. The improved scores in these metrics for PTAMT-enhanced models reflect diversified word usage and a departure from the simplified and often repetitive language characteristic of synthetic data-driven translations, thereby diminishing the hallmarks of machine translationese.

Besides lexical level, we have also observed a syntax-semantics phenomenon uniquely captured by the PTAMT-enhanced models. A translation sample (see Table 7 in Appendix C) is illustrated where the CA source sentence contains three elements conveying negative meaning, whereas the ZH reference exhibits only one negative marker. This is because CA is a negative concord language, where multiple negative markers do not cancel but affirm one another to intensify the negation, and thus combine into a single negation (Espinal et al., 2016; Tubau et al., 2023). By contrast, ZH is a language without negative concord, meaning that negative markers spell out one another and thus two negatives resolve to a positive (Yang, 2011). Therefore, the triple negatives in the CA source sentence actually resolve to a single negation. When translating the the CA source sentence to ZH, only one negative needs to be retained. In our experiments, only the PTAMT models accurately captured this linguistic phenomenon, while other models erroneously included two negatives in the ZH translation, resulting in a completely opposite meaning. Surprisingly, after fine-tuning the from-scratch trained PTAMT-enhanced model on the original parallel corpus, this understanding was

lost. Conversely, the M2M-100-418M PTAMT-enhanced model gained this understanding after the second-step fine-tuning.

Finally, we have also noted the substantial impact of incorporating a small quantity of authentic parallel data in the desired language pair (dataset *CA-ZH-WIKI*). When training or fine-tuning on pivot-translated data alongside a modest amount of original CA-ZH parallel corpus, there is a marked increase in spBLEU and COMET scores, compared to using pivot-translated data alone. Moreover, fine-tuning successively the from-scratch trained PTAMT-enhanced model and the M2M-100-418M PTAMT-enhanced model on a small portion of original CA-ZH parallel data also resulted in a notable enhancement in spBLEU and COMET scores for both models. Taken together, these findings are likely to imply the significance of authentic parallel data in the target language pair(s) in improving the performance of NMT systems.

6 Conclusions and Future Work

In this work, our comprehensive experimental evaluation of from-scratch trained and M2M-100-418M pre-trained models for the CA→ZH translation task has highlighted the efficacy of low-resource NMT techniques. Significantly, these experiments have confirmed the substantial benefits of these methods on translation quality and lexical diversity, with our novel PTAMT method emerging as a key innovator in addressing the challenges inherent in translating the low-resource language pair.

The PTAMT method, with its ability to effectively utilize pseudo-parallel and authentic parallel data, significantly mitigates the influence of source-side machine translationese and enhances the model’s capability to produce translations that are not only accurate but also linguistically rich and varied. This approach not only broadens the lexical range and usage contexts available to the model but also ensures a nuanced understanding and replication of natural language subtleties, as evidenced by the improved lexical diversity metrics and the accurate handling of complex linguistic phenomena such as negative concord.

Moreover, our findings seem to imply the critical role of integrating authentic data in the desired language pair(s) into the training or fine-tuning process, demonstrating that even a small amount

of authentic parallel data can substantially elevate the performance of NMT systems. This insight emphasizes the importance of combining pseudo-parallel and authentic inputs to achieve the best possible translation outcomes, particularly in the context of low-resource language pairs.

While our study marks progress in NMT for the low-resource CA-ZH pair, it also unveils areas which deserve further exploration. The potential domain alignment between our test set *FLORES-200* and the original *CA-ZH-WIKI* parallel data raises questions about how the inclusion of a modest amount of authentic parallel data in the target language pair(s) influences translation outcomes. This becomes especially relevant when considering the potential for domain-specific biases to affect the evaluation of NMT systems. Moreover, the observed discrepancies in how from-scratch trained PTAMT-enhanced models and M2M-100-418M PTAMT-enhanced models handle linguistic complexities such as negative concord—both before and after additional fine-tuning—suggest underlying differences in model learning dynamics that deserve closer scrutiny. These divergent model responses highlight the need for a nuanced understanding of how different training approaches impact the NMT models’ ability to grasp and accurately render complex linguistic structures.

To navigate these uncertainties and expand upon our findings, we propose several avenues for future research: firstly, building novel and diversified test sets to quantify and generalize the influence of authentic parallel data in the target language pair(s) on model performance; secondly, exploring the models’ internal representations and additional fine-tuning processes to pinpoint factors contributing to their distinct responses to linguistic complexities such as negative concord; thirdly, expanding our investigation to include more low-resource language pairs to enable a comprehensive evaluation of the PTAMT method’s applicability across diverse linguistic contexts.

7 Acknowledgements

This work was partly funded by the China Scholarship Council (CSC), to whom we express our sincere gratitude. We are also grateful to the anonymous reviewers of EAMT 2024, whose insightful comments significantly enhanced the quality and presentation of this paper.

References

- Adelani, David, Jesujoba Alabi, Angela Fan, Julia Kreutzer, Xiaoyu Shen, Machel Reid, Dana Ruitter, Dietrich Klakow, Peter Nabende, Ernie Chang, Tajuddeen Gwadabe, Freshia Sackey, Bonaventure F. P. Dossou, Chris Emezue, Colin Leong, Michael Beukman, Shamsuddeen Muhammad, Guyo Jarso, Oreen Yousuf, Andre Niyongabo Rubungo, Gilles Hacheme, Eric Peter Wairagala, Muhammad Umair Nasir, Benjamin Ajibade, Tunde Ajayi, Yvonne Gitau, Jade Abbott, Mohamed Ahmed, Millicent Ochieng, Anuoluwapo Aremu, Perez Ogayo, Jonathan Mukiiibi, Fatoumata Ouoba Kabore, Godson Kalipe, Derguene Mbaye, Allahsera Auguste Tapo, Victoire Memdjokam Koagne, Edwin Munkoh-Buabeng, Valencia Wagner, Idris Abdullumin, Ayodele Awokoya, Happy Buzaaba, Blessing Sibanda, Andiswa Bukula, and Sam Manthalu. 2022. A few thousand translations go a long way! leveraging pre-trained models for African news translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3053–3070, Seattle, United States. Association for Computational Linguistics.
- Araabi, Ali and Christof Monz. 2020. Optimizing transformer for low-resource neural machine translation.
- Arivazhagan, Naveen, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George Foster, Colin Cherry, Wolfgang Macherey, Zhifeng Chen, and Yonghui Wu. 2019. Massively multilingual neural machine translation in the wild: Findings and challenges.
- Chae, Jieun and Ani Nenkova. 2009. Predicting the fluency of text with shallow structural features: Case studies of machine translation and human-written text. In Lascarides, Alex, Claire Gardent, and Joakim Nivre, editors, *Proceedings of the 12th Conference of the European Chapter of the ACL (EACL 2009)*, pages 139–147, Athens, Greece, March. Association for Computational Linguistics.
- Chen, Yun, Yang Liu, Yong Cheng, and Victor O. K. Li. 2017. A teacher-student framework for zero-resource neural machine translation.
- Cohn, Trevor and Mirella Lapata. 2007. Machine translation by triangulation: Making effective use of multi-parallel corpora. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 728–735, Prague, Czech Republic. Association for Computational Linguistics.
- Costa-jussà, Marta Ruiz, José A. R. Fonollosa, José B. Mariño, Marc Poch, and Mireia Farrús. 2014. A large Spanish-Catalan parallel corpus release for machine translation. *Comput. Informatics*, 33:907–920.
- Costa-Jussà, Marta R., Noé Casas, Carlos Escolano, and José A. R. Fonollosa. 2019. Chinese-catalan: A neural machine translation approach based on pivoting and attention mechanisms. *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*, 18(4).
- Currey, Anna and Kenneth Heafield. 2019. Zero-resource neural machine translation with monolingual pivot data. In Birch, Alexandra, Andrew Finch, Hiroaki Hayashi, Ioannis Konstas, Thang Luong, Graham Neubig, Yusuke Oda, and Katsuhito Sudoh, editors, *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 99–107, Hong Kong, November. Association for Computational Linguistics.
- Dutta Chowdhury, Koel, Richa Jalota, Cristina España-Bonet, and Josef Genabith. 2022. Towards debiasing translation artifacts. In Carpuat, Marine, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3983–3991, Seattle, United States, July. Association for Computational Linguistics.
- El Kholly, Ahmed, Nizar Habash, Gregor Leusch, Evgeny Matusov, and Hassan Sawaf. 2013. Language independent connectivity strength features for phrase pivot statistical machine translation. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 412–418, Sofia, Bulgaria. Association for Computational Linguistics.
- El-Kishky, Ahmed, Vishrav Chaudhary, Francisco Guzmán, and Philipp Koehn. 2020. CCAIghned: A massive collection of cross-lingual web-document pairs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5960–5969, Online. Association for Computational Linguistics.
- Espinal, M. Teresa, Susagna Tubau, Joan Borràs-Comes, and Pilar Prieto, 2016. *Double Negation in Catalan and Spanish. Interaction Between Syntax and Prosody*, pages 145–176. Springer International Publishing, Cham.
- Fan, Angela, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Man-deep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation.
- Firat, Orhan, KyungHyun Cho, and Yoshua Bengio. 2016. Multi-way, multilingual neural machine translation with a shared attention mechanism. *CoRR*, abs/1601.01073.
- Freitag, Markus, Ricardo Rei, Nitika Mathur, Chikui Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022a. Results of WMT22 metrics shared

- task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Freitag, Markus, David Vilar, David Grangier, Colin Cherry, and George Foster. 2022b. A natural diet: Towards improving naturalness of machine translation output. In Muresan, Smaranda, Preslav Nakov, and Aline Villavicencio, editors, *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3340–3353, Dublin, Ireland, May. Association for Computational Linguistics.
- Goyal, Naman, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzman, and Angela Fan. 2021. The flores-101 evaluation benchmark for low-resource and multilingual machine translation.
- Gu, Jiatao, Hany Hassan, Jacob Devlin, and Victor O.K. Li. 2018. Universal neural machine translation for extremely low resource languages. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 344–354, New Orleans, Louisiana. Association for Computational Linguistics.
- Ha, Thanh-Le, Jan Niehues, and Alexander Waibel. 2016. Toward multilingual neural machine translation with universal encoder and decoder.
- Hassan Awadalla, Hany, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, Will Lewis, Mu Li, Shujie Liu, Tie-Yan Liu, Renqian Luo, Arul Menezes, Tao Qin, Frank Seide, Xu Tan, Fei Tian, Lijun Wu, Shuangzhi Wu, Yingce Xia, Dongdong Zhang, Zhirui Zhang, and Ming Zhou. 2018. Achieving human parity on automatic chinese to english news translation.
- Ilisei, Iustina, Diana Inkpen, Gloria Corpas Pastor, and Ruslan Mitkov. 2010. Identification of translationese: A machine learning approach. In Gelbukh, Alexander, editor, *Computational Linguistics and Intelligent Text Processing*, pages 503–511, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Jalota, Rricha, Koel Chowdhury, Cristina España-Bonet, and Josef van Genabith. 2023. Translating away translationese without parallel data. In Bouamor, Houda, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7086–7100, Singapore, December. Association for Computational Linguistics.
- Johnson, Melvin, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google’s multilingual neural machine translation system: Enabling zero-shot translation.
- Kreutzer, Julia, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, Monang Setyawan, Supheakmongkol Sarin, Sokhar Samb, Benoît Sagot, Clara Rivera, Annette Rios, Isabel Papadimitriou, Salomey Osei, Pedro Ortiz Suarez, Iro Orife, Kelechi Ogueji, Andre Niyongabo Rubungo, Toan Q. Nguyen, Mathias Müller, André Müller, Shamsuddeen Hassan Muhammad, Nanda Muhammad, Ayanda Mnyakeni, Jamshidbek Mirzakhlov, Tapiwanashe Matangira, Colin Leong, Nze Lawson, Sneha Kudugunta, Yacine Jernite, Mathias Jenny, Orhan Firat, Bonaventure F. P. Dossou, Sakhile Dlamini, Nisansa de Silva, Sakine Çabuk Ballı, Stella Biderman, Alessia Battisti, Ahmed Baruwa, Ankur Bapna, Pallavi Baljekar, Israel Abebe Azime, Ayodele Awokoya, Duygu Ataman, Orevaoghene Ahia, Oghenefego Ahia, Sweta Agrawal, and Mofetoluwa Adeyemi. 2022. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72.
- Lakew, Surafel Melaku, Quintino Francesco Lotito, Matteo Negri, Marco Turchi, and Marcello Federico. 2018. Improving zero-shot translation of low-resource languages. *ArXiv*, abs/1811.01389.
- Liu, Zixuan. 2022. Improving chinese-catalan machine translation with wikipedia parallel corpus. Master’s thesis, Universitat Pompeu Fabra, Barcelona.
- Lu, Yichao, Phillip Keung, Faisal Ladhak, Vikas Bhardwaj, Shaonan Zhang, and Jason Sun. 2018. A neural interlingua for multilingual machine translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 84–92, Brussels, Belgium. Association for Computational Linguistics.
- Neubig, Graham and Junjie Hu. 2018. Rapid adaptation of neural machine translation to new languages. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 875–880, Brussels, Belgium. Association for Computational Linguistics.
- Popel, Martin, Markéta Tomková, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar, and Zdeněk Zábokrtský. 2020. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. *Nature Communications*, 11:4381.
- Rapp, Reinhard. 2021. Similar language translation for Catalan, Portuguese and Spanish using Marian NMT. In *Proceedings of the Sixth Conference on Machine Translation*, pages 292–298, Online. Association for Computational Linguistics.

- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.
- Riley, Parker, Isaac Caswell, Markus Freitag, and David Grangier. 2020. Translationese as a language in “multilingual” nmt.
- Schwenk, Holger, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2019. Wiki-matrix: Mining 135m parallel sentences in 1620 language pairs from wikipedia.
- Schwenk, Holger, Guillaume Wenzek, Sergey Edunov, Edouard Grave, and Armand Joulin. 2020. Cmatrix: Mining billions of high-quality parallel sentences on the web.
- Tan, Xu, Yi Ren, Di He, Tao Qin, Zhou Zhao, and Tie-Yan Liu. 2019. Multilingual neural machine translation with knowledge distillation.
- Team, NLLB, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.
- Thillainathan, Sarubi, Surangika Ranathunga, and Sanath Jayasena. 2021. Fine-tuning self-supervised multilingual sequence-to-sequence models for extremely low-resource nmt. In *2021 Moratuwa Engineering Research Conference (MERCon)*, pages 432–437.
- Tubau, Susagna, Urtzo Etxeberria, and M. Teresa Espinal. 2023. A new approach to negative concord: Catalan as a case in point. *Journal of Linguistics*, page 1–33.
- Vanmassenhove, Eva, Dimitar Shterionov, and Matthew Gwilliam. 2021. Machine translationese: Effects of algorithmic bias on linguistic complexity in machine translation. In Merlo, Paola, Jorg Tiedemann, and Reut Tsarfaty, editors, *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2203–2213, Online, April. Association for Computational Linguistics.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.
- Wang, Xinyi, Yulia Tsvetkov, and Graham Neubig. 2020. Balancing training for multilingual neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8526–8537, Online. Association for Computational Linguistics.
- Wang, Rui, Xu Tan, Renqian Luo, Tao Qin, and Tie-Yan Liu. 2021. A survey on low-resource neural machine translation.
- Wein, Shira and Nathan Schneider. 2024. Lost in translationese? reducing translation effect using abstract meaning representation.
- Wu, Hua and Haifeng Wang. 2007. Pivot language approach for phrase-based statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 856–863, Prague, Czech Republic. Association for Computational Linguistics.
- Yang, Huiling. 2011. Is Chinese a negative concord language? In *proceedings of the 23rd North American Conference on Chinese Linguistics (NACCL-23)*, pages 208–223, Eugene, United States), December. University of Oregon Press.
- Zheng, Hao, Yong Cheng, and Yang Liu. 2017. Maximum expected likelihood estimation for zero-resource neural machine translation. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence, IJCAI’17*, page 4251–4257. AAAI Press.
- Zhou, Chenuye. 2022. Building a catalan-chinese parallel corpus from Wikipedia for use in machine translation. Master’s thesis, Universitat Pompeu Fabra, Barcelona.
- Ziemski, Michał, Marcin Junczys-Dowmunt, and Bruno Pouliquen. 2016. The United Nations parallel corpus v1.0. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, Portorož, Slovenia. European Language Resources Association (ELRA).
- Zoph, Barret, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation.

A Appendix: Hyper-parameters Configuration

Hyper-parameters	Value
adam betas	0.9, 0.98
learning rate	0.0005
warmup initial learning rate	1.00E-07
label smoothing	0.1
dorpout	0.2
weight decay	0.0001
batch size (in tokens)	4096 (\times 8 steps)
gradients accumulation	8

Table 4: hyper-parameters of the neural models. Note that the same set of hyper-parameters was used for all experiments except that the batch size for the M2M-100-418M models was 2048 tokens (\times 16 steps) due to GPU memory limit.

Hyperparameter	Value
adam betas	0.9, 0.98
learning rate	0.0005
warmup initial learning rate	1.00E-07
label smoothing	0.5
dorpout	0.3
activation dropout	0.3
enc/dec layerDrop	0.0/0.1
weight decay	0.0001
batch size (in tokens)	4096 (\times 8 steps)
gradients accumulation	8

Table 5: Optimal Transformer hyper-parameters settings for 40k datasets.

B Appendix: Training and Fine-tuning Epochs

System	Methods	Epoch
Vanilla Baseline	-	155
M2M-100-418M Models	multilingual training (strong baseline)	-
	fine-tuning	8
	fine-tuning & pivot (without original data)	3
	fine-tuning & pivot (with original data)	4
	PTAMT	4
	PTAMT & 2nd fine-tuning	10
	From-scratch Trained Models	pivot (without original data)
pivot (with original data)		11
PTAMT		13
PTAMT & 2nd fine-tuning		10

Table 6: Training and fine-tuning epochs for the baseline, pre-trained models and from-scratch models. Note that the epoch of the pre-trained M2M-100-418M model was not reported publicly.

C Appendix: Translation Sample

Model	Sentence
Source(CA)	Adverteix que <u>no</u> hi ha <u>ningú</u> que pugui garantir que <u>cap</u> acció a l'Iraq en aquest moment aconseguixi aturar la guerra sectària, la violència creixent o una deriva caòtica.
Reference(ZH)	报告警告称,没有人能保证目前在伊拉克采取的任何行动能够阻止宗派战争、不断增长的暴力或走向混乱。
<i>Meaning in English</i>	It warns <u>no</u> one can guarantee that any action in Iraq at this point will stop sectarian warfare, growing violence, or a slide toward chaos.
Vanilla baseline	亚当斯表示,没有任何人确保伊拉克的行动,并阻止伊拉克战争、暴力行为、暴力行为或暴力行为。
<i>Meaning in English</i>	Adam shows, <u>no</u> one guarantees the action in Iraq and stops the war in Iraq, violence, violence, or violence.
M2M-100-418M (Strong baseline)	他警告说,没有人能保证伊拉克目前没有任何行动能阻止种族战争、暴力加剧或混乱的流动。
<i>Meaning in English</i>	He warns <u>no</u> one can guarantee that <u>no</u> action in Iraq at this point will stop ethnic warfare, growing violence, or the flow of chaos.
M2M-100-418M + finetuning	他警告,没有人能保证伊拉克在这一时刻不会采取任何行动来阻止种族战争、日益暴力或混乱的发生。
<i>Meaning in English</i>	He warns <u>no</u> one can guarantee that Iraq is <u>not</u> taking any action at this point to stop ethnic warfare, growing violence, or the occurrence of chaos.
M2M-100-418M + fine-tuning & pivot (without original data)	他警告说,没有人能够保证目前在伊拉克采取的任何行动都不会停止派别战争、暴力升级或混乱。
<i>Meaning in English</i>	He warns <u>no</u> one can guarantee that any action in Iraq at this point will <u>not</u> stop sectarian warfare, growing violence, or chaos.
M2M-100-418M + fine-tuning & pivot (with original data)	他警告说,没有人能够保证目前在伊拉克采取的任何行动都不会阻止派别战争、暴力升级或混乱。
<i>Meaning in English</i>	He warns <u>no</u> one can guarantee that any action in Iraq at this point will <u>not</u> stop sectarian warfare, growing violence, or chaos.
M2M-100-418M + PTAMT	他警告说,没有人能够保证,目前在伊拉克采取的任何行动都不会成功地制止教派战争、不断升级的暴力或混乱的倾向。
<i>Meaning in English</i>	He warns <u>no</u> one can guarantee that any action in Iraq at this point will <u>not</u> successfully stop sectarian warfare, growing violence, or tendency towards chaos.

Model	Sentence
M2M-100-418M + PTAMT & 2nd-fine-tuning <i>Meaning in English</i>	他警告说,没有人能够保证在伊拉克的任何行动能阻止教派战争、不断升级的暴力或混乱的倾向。 He warns <u>no</u> one can guarantee that any action in Iraq will stop sectarian warfare, growing violence, or tendency towards chaos.
From-scratch + pivot translation (without original data) <i>Meaning in English</i>	他指出,没有人能够保证伊拉克目前的任何行动都不会导致教派战争、日益严重的暴力或混乱。 He points out <u>no</u> one can guarantee that any action in Iraq at this point will <u>not</u> lead to sectarian warfare, growing violence, or chaos.
From-scratch + pivot translation (with original data) <i>Meaning in English</i>	他警告说,没有人能够保证,伊拉克目前的任何行动都不会阻止教派战争、不断升级的暴力或混乱的漂流。 He warns <u>no</u> one can guarantee that any action in Iraq at this point will <u>not</u> stop sectarian warfare, growing violence, or the flow of chaos.
From-scratch + PTAMT <i>Meaning in English</i>	他警告说,没有人能够确保目前在伊拉克的任何行动能够制止教派战争、不断升级的暴力或混乱。 He warns <u>no</u> one can guarantee that any action in Iraq at this point will stop sectarian warfare, growing violence, or chaos.
From-scratch + PTAMT & 2nd-fine-tuning <i>Meaning in English</i>	他警告说,没有能保证此时在伊拉克的任何行动都不会阻止教派战争、不断升级的暴力或混乱的漂流。 He warns <u>no</u> one can guarantee any action in Iraq at this point will <u>not</u> stop sectarian warfare, growing violence, or the flow of chaos.

Table 7: Translation sample for baseline, fine-tuned, and from-scratch trained models. Note that the underlined elements in the table are words or structural elements that cause negation.

A Case Study on Context-Aware Neural Machine Translation with Multi-Task Learning

Ramakrishna Appicharla¹, Baban Gain¹, Santanu Pal², Asif Ekbal³, Pushpak Bhattacharyya⁴

¹Department of Computer Science and Engineering, Indian Institute of Technology Patna, India

²Wipro AI, Lab45, London, UK

³School of AI and Data Science, Indian Institute of Technology Jodhpur, India

⁴Department of Computer Science and Engineering, Indian Institute of Technology Bombay, India

{ramakrishnaappicharla, gainbaban, santanu.pal.ju,
asif.ekbal, pushpakbh} @gmail.com

Abstract

In document-level neural machine translation (DocNMT), multi-encoder approaches are common in encoding context and source sentences. Recent studies (Li et al., 2020) have shown that the context encoder generates noise and makes the model robust to the choice of context. This paper further investigates this observation by explicitly modelling context encoding through multi-task learning (MTL) to make the model sensitive to the choice of context. We conduct experiments on cascade MTL architecture, which consists of one encoder and two decoders. Generation of the source from the context is considered an auxiliary task, and generation of the target from the source is the main task. We experimented with German–English language pairs on News, TED, and Europarl corpora. Evaluation results show that the proposed MTL approach performs better than concatenation-based and multi-encoder DocNMT models in low-resource settings and is sensitive to the choice of context. However, we observe that the MTL models are failing to generate the source from the context. These observations align with the previous studies, and this might suggest that the available document-level parallel corpora are not context-aware, and a robust sentence-level model can outperform the context-aware models.

1 Introduction

Context-aware neural machine translation gained much attention due to the ability to incorporate context, which helps in producing more consistent translations than sentence-level models (Maruf and Haffari, 2018; Zhang et al., 2018; Bawden et al., 2018; Agrawal et al., 2018; Voita et al., 2019; Huo et al., 2020; Li et al., 2020; Donato et al., 2021). There are mainly two approaches to incorporating context. The first one is to create a context-aware input sentence by concatenating context and current input sentence (Tiedemann and Scherrer, 2017; Agrawal et al., 2018; Junczys-Dowmunt, 2019; Zhang et al., 2020b) and using it as the input to the encoder. The second approach uses an additional context-aware component to encode the source or target context (Zhang et al., 2018; Voita et al., 2018; Kim et al., 2019; Ma et al., 2020) and the entire model is jointly optimized. Typically, the current sentence’s neighbouring sentences (either previous or next) are used as the context.

The context-aware models are trained to maximize the log-likelihood of the target sentence given the source sentence and context. Most of the existing works on DocNMT (Zhang et al., 2018; Maruf and Haffari, 2018; Voita et al., 2019; Li et al., 2020) focus on encoding the context through context-specific encoders. Recent studies (Li et al., 2020) show that, in the multi-encoder DocNMT models, the performance improvement is not due to specific context encoding but rather the context-encoder acts like a noise generator, which, in turn, improves the robustness of the model. In this work, we explore whether the context encoding can be modelled explicitly through multi-task learning (MTL) (Luong et al., 2015). Specifically, we aim to study the effectiveness of the MTL framework

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

for DocNMT rather than proposing a state-of-the-art system. The availability of document-level corpora is less compared to sentence-level corpora. Previous works (Junczys-Dowmunt, 2019) use the sentence-level corpora to warm-start the document-level model, which can be further tuned with the existing limited amount of document-level data. However, in this work, we focus only on improving the performance of DocNMT models with available document-level corpora. We consider the source reconstruction from the context as the auxiliary task and the target translation from the source as the main task. We conduct experiments on cascade MTL (Anastasopoulos and Chiang, 2018; Zhou et al., 2019) architecture. The cascade MTL architecture comprises one encoder and two decoders (Figure 1). The intermediate (first) decoder attends over the output of the encoder, and the final (second) decoder attends over the output of the intermediate decoder. The input consists of $\langle c_x, x, y \rangle$ triplets, where c_x , x and y represents the context, source, and target sentences, respectively. The model is trained to optimize both translation and reconstruction objectives jointly. We also train two baseline models as contrastive models, namely sentence-level vanilla baseline and single encoder-decoder model, by concatenating the context and source (Tiedemann and Scherrer, 2017; Agrawal et al., 2018; Junczys-Dowmunt, 2019). We additionally train multi-encoder single-decoder models (Li et al., 2020) to study how context affects the DocNMT models. We conduct experiments on German-English direction with three different types of contexts (*viz.* previous two source sentences, previous two target sentences, and previous-next source sentences) on News-commentary v14 and TED corpora. We report BLEU (Papineni et al., 2002) calculated with sacreBLEU (Post, 2018) and APT (accuracy of pronoun translation) (Miculicich Werlen and Popescu-Belis, 2017) scores.

To summarize, the specific attributes of our current work are as follows:

- We explore whether the MTL approach can improve the performance of context-aware NMT by introducing additional training objectives along with the main translation objective.
- We propose an MTL approach where the reconstruction of the source sentence given the

context is used as an auxiliary task and the translation of the target sentence from the source sentence as the main task, jointly optimized during the training.

- The results show that in the MTL approach, the context encoder generates noise similar to the multi-encoder approach (Li et al., 2020), which makes the model robust to the choice of the context.

2 Related Work

Previous studies have proposed various document-level NMT models and achieved great success. The main goal of these approaches is to efficiently model context representation, which can lead to better translation quality. Towards this goal to represent context, Tiedemann and Scherrer (2017) concatenate consecutive sentences and use them as input to the single-encoder-based DocNMT model. Agrawal et al. (2018) conducted experiments on varying neighbouring contexts and concatenated with the current sentence as input to their model. With these similar trends, Junczys-Dowmunt (2019) conducted experiments considering the entire document as context. Further progress on context representation in DocNMT, Zhang et al. (2018) and Voita et al. (2018) proposed transformer-based multi-encoder NMT models where the additional encoder is used to encode the context. While Miculicich et al. (2018) proposed a hierarchical attention network to encode the context, a more recent approach Kang et al. (2020) proposed a reinforcement learning-based dynamic context selection module for DocNMT. Kim et al. (2019) and Li et al. (2020) conducted experiments on multi-encoder DocNMT models and reported that the performance improvement is not due to context encoding; instead, the context encoder acts as a noise generator, which improves the robustness of the DocNMT model. Junczys-Dowmunt (2019) conducted experiments on a single encoder model with masked language model objective (Devlin et al., 2019) to incorporate document-level monolingual source-side data. Since the multi-encoder models are trained to optimize the translation objective only, it might be possible for the model to pay less attention to the context, and Li et al. (2020) report the same.

MTL strategies in NMT trained on other auxiliary tasks along with the main translation task (Lu-

ong et al., 2015; Dong et al., 2015; Zareemoodi et al., 2018; Wang et al., 2020; Yang et al., 2020) achieved significant improvements in translation quality so far. The other auxiliary tasks include autoencoding (Luong et al., 2015), denoising autoencoding (Wang et al., 2020), parsing and named entity recognition (Zareemoodi and Haffari, 2018; Zareemoodi et al., 2018). Zhou et al. (2019) proposed a cascade MTL network to improve the robustness of the NMT model. They considered denoising the noisy text as an auxiliary task and the translation as the main task. They achieved a significant BLEU score improvement (up to 7.1 BLEU) on the WMT robustness shared task on the French-English dataset.

However, most multi-task models are proposed only for sentence-level NMT models. Multi-task learning is relatively unexplored in context-aware NMT settings. Wang et al. (2021) proposed an MTL framework for dialogue translation tasks that jointly correct the sentences having issues such as pronoun dropping, punctuation dropping, and typos and translate them into the target language. Liang et al. (2022) proposed a three-stage training framework for the neural chat translation task. The model is trained on auxiliary tasks such as monolingual cross-lingual response generation tasks to generate coherent translation and the next utterance discrimination task. Lei et al. (2022) proposed an MTL system to force the model to attend over relevant cohesion devices while translating the current sentence. In this work, we propose a multi-task learning objective, i.e., reconstruction of source sentences given the source context in a cascade multi-task learning setting to study the effect of context in document-level NMT systems.

3 Methodology

3.1 Problem Statement

Our document-level NMT is based on a cascade MTL framework to force the model to consider the context while generating translation. Given a source sentence x and context c_x , the translation probability of the target sentence y in the DocNMT setting is calculated as in Equation 1.

$$p(y) = p(y|x, c_x) \times p(x, c_x) \quad (1)$$

We consider $p(x, c_x)$ as the auxiliary task of source (x) reconstruction from c_x (as $p(x|c_x)$)¹,

¹Since the joint probability of $p(x, c_x)$ can be calculated as

calculated as in Equation 2.

$$p(x, c_x) = p(x|c_x) \times p(c_x) \quad (2)$$

The training data D consists of triplets $\langle c_x, x, y \rangle$. Given the parameters of the model θ , the translation (Equation 1) and reconstruction (Equation 2) objectives can be modeled as Equation 3 and Equation 4.

$$p(y|x, c_x; \theta) = \prod_{t=1}^T p(y_t|x, c_x, y_{<t}; \theta) \quad (3)$$

$$p(x|c_x; \theta) = \prod_{s=1}^S p(x_s|c_x, x_{<s}; \theta) \quad (4)$$

where, S, Z, T denote the lengths of x, c_x, y respectively and $x_{<s}, c_{x<z}, y_{<t}$ denote partially generated sequences.

Given translation objective $p(y|x, c_x)$ and reconstruction objective $p(x|c_x)$, the model is jointly trained and optimized the loss, \mathcal{L} using parameter θ (cf. Equation 5); where α is a hyper-parameter used to control the loss. We set α to 0.5.

$$\mathcal{L} = \alpha * \log p(y|x, c_x; \theta) + (1 - \alpha) * \log p(x|c_x; \theta) \quad (5)$$

We hypothesize that forcing the model to learn reconstruction and translation objectives jointly will enable the model to encode the context effectively. The output of the reconstruction task can verify this during testing. If the context encoder generates noise, then the model might be unable to reconstruct the source and vice-versa.

3.2 Cascade Multi-Task Learning Transformer

The cascade multi-task learning architecture (Zhou et al., 2019) (Figure 1) consists of one encoder and two decoders based on the transformer (Vaswani et al., 2017) architecture. The model takes three inputs: *Source*: Current source sentence, *Context*: Context of the current source sentence, and *Target*: Current target sentence. The input to the encoder is context, and the input to the intermediate decoder is the source. The intermediate decoder is trained

$p(c_x|x) \times p(x)$, we also explored this setting. We observed that the performance of the model is poor in this setting compared to the other setting. More details can be found in Appendix A.1.

to reconstruct the source given context by attending to the output of the encoder. The final decoder attends over the output of the intermediate decoder. In the non-MTL setting, the model is trained only on the translation objective (output of the final decoder), and the intermediate decoder is not trained with the reconstruction objective.

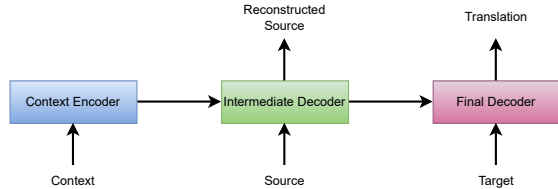


Figure 1: The overview of our MTL architecture. The input to the model is a triplet. The triplet consist of (*Context*, *Source*, *Target*). The Intermediate Decoder is trained to reconstruct the *Source* given *Context*, and the Final Decoder is trained to translate the *Source*. Here, *Source*: Current source sentence, *Context*: Context for the current source sentence, and *Target*: Translation of current source sentence. None of the layers are shared.

3.3 Context Selection

We conduct experiments on different settings of the source context. The term “source context” is defined as considering related or dependent sentences directly related to the input sentence. Based on the findings of Zhang et al. (2018), we select two sentences as context and concatenate them with a special token ‘<break>’ (Junczys-Dowmunt, 2019). For a given input source sentence (x_i) and target sentence (y_i), contexts selected for the experiments are:

- Previous-2 Source (**P@2-SRC**): Two previous source sentences (x_{i-2}, x_{i-1})
- Previous-2 Target (**P@2-TGT**): Two previous target sentences (y_{i-2}, y_{i-1})
- Previous-Next Source (**P-N-SRC**): Previous and next source sentences (x_{i-1}, x_{i+1})

4 Experiment Setup

We train our models with the proposed cascade MTL approach. The model is trained on $\langle c_x, x, y \rangle$ triplet to jointly optimize both translation and source reconstruction objectives (Figure: 1). We also train three other contrastive models to show the effect of context in the MTL setting.

Vanilla-Sent: A vanilla sentence-level baseline model is trained without context on a single encoder-decoder network.

Concat-Context: This model is trained on a single encoder-decoder network where context is concatenated with the source (Tiedemann and Scherrer, 2017; Agrawal et al., 2018; Junczys-Dowmunt, 2019) and fed to the encoder as input. In this setting, sentences within the context are concatenated with a unique token, ‘<break>’. The context and the source are concatenated with another special symbol, ‘<concat>’. The special symbol helps the model to distinguish between context and source sentences.

Inside-Context: We re-implemented the ‘Inside-Context’ model proposed by Li et al. (2020), a multi-encoder approach. This model consists of two encoders and one decoder. The decoder is modified to include two cross-attention layers to attend over the outputs of both encoders before passing through the position-wise feed-forward layer (Vaswani et al., 2017).

4.1 Data Statistics

We conduct experiments on WMT news-commentary, IWSLT’17 TED, and Europarl-v7 German-English corpora. For the WMT news-commentary, we use news-commentary v14 (Barrault et al., 2019)² as the train set, newstest2017 as the validation set, and newstest2018 as the test set. For IWSLT’17 TED and Europarl-v7 corpora, we follow the train, validation, and test set splits mentioned in (Maruf et al., 2019)³. All models are trained on German to English. Table 1 shows data statistics of the train, validation, and test sets.

Data	# Sent	# Doc
News	329,000/3,004/2,998	8,462/130/122
TED	206,112/8,967/2,271	1,698/93/23
Europarl	1,666,904/3,587/5,134	117,855/240/360

Table 1: Data statistics for our experiments. # Sent, # Doc represent the number of sentences and documents, respectively. The numbers are shown in the Train/Validation/Test set order.

4.2 NMT Model Setups

We conduct all the experiments on transformer architecture (Vaswani et al., 2017). All the mod-

²<https://data.statmt.org/news-commentary/v14/training/>

³<https://github.com/sameenmaruf/selective-attn/tree/master/data>

Model	News		TED		Europarl	
	s-BLEU	d-BLEU	s-BLEU	d-BLEU	s-BLEU	d-BLEU
Vanilla-Sent	18.3	20.9	19.9	24.9	32.3	35.1
Concat-Context: P@2-SRC	18.0	20.5	17.3	22.4	32.5	35.4
Concat-Context: P-N-SRC	18.4	20.7	17.5	22.5	32.7	35.6
Concat-Context: P@2-TGT	14.7	17.2	15.3	20.4	36.4	39.1
MTL: P@2-SRC	19.1	21.7	20.2	24.8	29.5	32.6
MTL: P-N-SRC	20.1 [†]	22.5	20.3	25.2	32.5 [†]	35.3
MTL: P@2-TGT	19.2	21.7	20.7 [†]	25.4	28.2	31.6

Table 2: BLEU scores of Vanilla-Sent, Concat-Context, and proposed MTL DocNMT models trained with different source contexts for German to English direction on News-commentary v14, IWSLT-17 TED, and Europarl corpora. **s-BLEU** and **d-BLEU** represent sentence-level and document-level BLEU respectively. The best results are shown in bold. ‘†’ denotes the statistically significant results than Vanilla-Sent and Concat-Context models with $p < 0.05$.

els are implemented in PyTorch⁴. We use 6-layer encoder-decoder stacks with 8 attention heads. Positional token embedding sizes are set to 512, and the feed-forward layer consists of 2048 cells. Adam optimizer (Kingma and Ba, 2015) is used for training with a noam learning rate scheduler (Vaswani et al., 2017) with an initial learning rate of 0.2. We use warmup steps of 16,000 (Popel and Bojar, 2018), and dropout is set to 0.1. Due to the GPU memory restrictions, we use a mini-batch of 40 sentences for the models trained on News and TED corpora and 25 for the models trained on Europarl corpus. We create joint subword vocabularies of size 32k for each training corpus. We use the BPE (Sennrich et al., 2016) to create subword vocabularies with SentencePiece (Kudo and Richardson, 2018) implementation. We also learn the positional encoding of tokens (Devlin et al., 2019), and the maximum sequence length is set to 140 tokens for all models and 160 for *Concat-Context* models.

All the models are trained till convergence. We use the perplexity of the validation set as an early stopping criterion with the patience of 10 (Popel and Bojar, 2018). We report results on the best model checkpoint saved during the training. We perform beam search during inference with beam size 4 and length penalty of 0.6 (Wu et al., 2016). For DocNMT models, we use the same source context with which the models are trained. Since the input to the intermediate decoder (source sentence) is also given during the testing phase, the representation of the intermediate decoder can be calculated in parallel, similar to the training phase.

All the experiments are conducted on a single Nvidia GTX 2080ti GPU. The number of parameters and training time of the models is as follows:

Vanilla-Sent: 76M, 76.5 hours, *Concat-Context*: 76M, 81 hours, *Inside-Context*: 118M, 125 hours and proposed *MTL*: 130M, 160 hours. The parameters and training times are approximately the same for all the corpora.

5 Results and Analysis

This section discusses the results of the trained models and the context’s effect on Multi-Encoder and MTL settings. Table 2 shows the sentence-BLEU (s-BLEU) and document-BLEU (d-BLEU) (Liu et al., 2020; Bao et al., 2021) scores of the proposed multi-task learning model along with the *Vanilla-Sent* and *Concat-Context* models.

We report all models’ BLEU scores on German \rightarrow English direction, calculated with sacreBLEU (Post, 2018).

5.1 Results of MTL and Contrastive Models

We report the BLEU scores of the models on German \rightarrow English direction, calculated with sacreBLEU (Post, 2018)⁵. The proposed MTL model can outperform both *Vanilla-Sent* and *Concat-Context* models by achieving s-BLEU scores of 20.1 (*MTL: P-N-SRC*) and 20.7 (*MTL: P@2-TGT*) with an improvement of +1.8 and +0.8 BLEU improvement for News and TED corpora respectively. However, in the case of the Europarl data set, *Concat-Context* models outperform both *Vanilla-Sent* and *MTL* models. This shows that the *Concat-Context* model requires more data to perform well, unlike the MTL models, which can also work effectively in low-resource settings. We observe that the performance of the models is almost uniform across the three different context settings

⁴<https://pytorch.org/>

⁵sacreBLEU signature:“nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1”

with a maximum BLEU difference of +1.0 ($P\text{-}N\text{-}SRC$ vs. $P@2\text{-}SRC$) on News, +0.5 ($P@2\text{-}TGT$ vs. $P@2\text{-}SRC$) on TED and +4.3 ($P\text{-}N\text{-}SRC$ vs $P@2\text{-}TGT$) on Europarl corpora respectively.

We also report d-BLEU (document-level BLEU) scores (Liu et al., 2020; Bao et al., 2021) by converting each document into one single sequence (paragraph) by concatenating all sentences from that document and calculate BLEU scores on the resulting corpus. This results in slightly higher scores than the sentence level by matching n-grams over the whole document instead of at the sentence level. Table 2 also shows d-BLEU scores. Like s-BLEU scores, proposed MTL models achieve the best d-BLEU scores of 22.5 and 25.4 for News and TED corpora, respectively. We report the paired bootstrap resampling (Koehn, 2004) results, calculated with sacreBLEU (Post, 2018).

Model	News	TED	Europarl
MTL: P@2-SRC	1.3	1.4	4.9
MTL: P@2-TGT	1.2	1.6	3.9
MTL: P-N-SRC	1.3	1.5	3.1

Table 3: s-BLEU scores for the reconstruction objective of the MTL models on test set for News, TED, and Europarl corpora.

5.2 Analysis of Reconstruction Objective

We analyze the performance of the MTL model on the reconstruction objective on the test set to verify if the context encoder is generating noise. If the context encoder generates noise by the suboptimal encoding of context, the intermediate decoder will fail to reconstruct the source sentence from the context; otherwise, the intermediate decoder can reconstruct the source sentence to a similar extent as the final translated sentence. We perform greedy decoding on the intermediate decoder to generate the source from the context. Table 3 shows the BLEU scores of the reconstruction objective on the test set for News, TED, and Europarl corpora. The results show that the MTL models fail to reconstruct the source from the context. Based on this, we conclude that the context encoder cannot encode the context, leading to poor reconstruction performance of the models. However, we hypothesize that the model cannot reconstruct the source from the context because the corpora used to train context-aware models might not be context-aware. This observation aligns with the previous works

(Kim et al., 2019; Li et al., 2020), and with enough data, vanilla sentence-level NMT models can outperform the document-level NMT models.

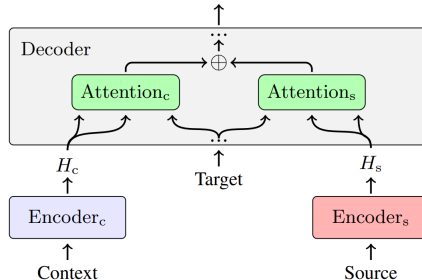


Figure 2: The overview of the Inside-Context model. The input to the model is a triplet consisting of ($Context$, $Source$, $Target$). The multi-head attention layer of the decoder is modified to attend to both the context encoders ($Encoder_c$) and the source encoder ($Encoder_s$).

Model	News	TED	Europarl
MTL: P@2-SRC	19.1	20.2	29.5
MTL: P-N-SRC	20.1 [†]	20.3	32.5 [†]
MTL: P@2-TGT	19.2	20.7 [†]	28.2
Inside-Context: P@2-SRC	18.8	19.6	33.2
Inside-Context: P-N-SRC	19.0	19.8	33.2
Inside-Context: P@2-TGT	18.3	20.4	33.6

Table 4: Comparison of s-BLEU scores of MTL and Inside-Context Multi-Encoder models. The best results are shown in bold. ‘[†]’ denotes the statistically significant results than Vanilla-Sent and Concat-Context models with $p < 0.05$.

5.3 MTL vs. Multi-Encoder Approach

We compare the proposed MTL approach to the existing Multi-Encoder approach to study how the model will perform in a single-task setting. Specifically, we compare our MTL approach (single-encoder multi-decoder network) with *Inside-Context* (Li et al., 2020) architecture. This model consists of two transformer encoders and one transformer decoder. Figure 2 shows the model’s architecture. The decoder is modified to attend to the outputs of both encoders. The model follows the transformer (Vaswani et al., 2017) architecture. An element-wise addition is performed on the outputs of both cross-attention layers before passing through layer-norm and position-wise feed-forward layers. Table 4 shows the s-BLEU scores of the MTL and Inside-Context models. We observe that the performance of multi-encoder models is similar to MTL models, with MTL models achieving +1.1 (P-N-SRC models), +0.3 (P@2-TGT models) BLEU points im-

provement over Inside-Context models for News and TED corpora respectively. In the case of Europarl, inside-context models achieve better performance than the MTL models, with the P@2-TGT model achieving +5.4 BLEU points improvement compared to the MTL model. Based on the results, we conclude that the MTL setting is more effective for low-resource scenarios.

Model	News	TED	Europarl
MTL: P@2-SRC	1.2 (-17.9)	0.8 (-19.4)	4.5 (-25.0)
MTL: P-N-SRC	1.2 (-18.9)	0.8 (-19.5)	4.0 (-28.5)
MTL: P@2-TGT	0.5 (-18.7)	0.3 (-20.4)	3.9 (-24.3)
Inside-Context: P@2-SRC	18.7 (-0.1)	19.4 (-0.2)	33.2 (0.0)
Inside-Context: P-N-SRC	18.9 (-0.1)	19.8 (0.0)	33.2 (0.0)
Inside-Context: P@2-TGT	18.3 (0.0)	20.3 (-0.1)	33.1 (-0.5)

Table 5: Comparison of s-BLEU scores of MTL models tested with random context. The difference in scores over the models trained with the selected context is shown inside the parentheses.

5.4 Effect of Context in MTL setting

Since the BLEU scores of our MTL models are almost the same for all three context settings, we check whether the MTL models are affected by the choice of context. To this end, we test the MTL models with random context. Here, random context denotes two randomly selected sentences from the entire corpus. Table 5 shows the results of MTL and Inside-Context models tested with random context. Results show that the MTL models fail to translate source sentences when the context is random. However, Inside-Context models are agnostic to context as models can translate well even if the context is random. Our findings in the case of multi-encoder models are in line with the findings of Li et al. (Li et al., 2020). Based on the results, we conclude that MTL models are sensitive to the choice of context. Section A.1.1 describes a similar experiment where the MTL models are tested with random context. However, the architecture used in the main experiments differs slightly from the one used in the preliminary investigation. We observe that feeding the Intermediate Decoder output to the Final Decoder makes the model sensitive to the choice of context (cf. Figure 1 and Figure 3 in the Appendix A.1). We hypothesize that a weighted combination of the Context Encoder output and Intermediate Decoder output is desired as it performs slightly better than the model used in the main experimental setup. However, it also makes the model agnostic to the choice

of context. We plan to explore this behaviour in detail in our future work.

Model	News	TED	Europarl
MTL: P@2-SRC	13.7 (+12.5)	11.2 (+10.4)	22.3 (+17.8)
MTL: P-N-SRC	14.5 (+13.3)	11.3 (+10.5)	19.7 (+15.7)
Inside-Context: P@2-SRC	18.7 (0.0)	19.6 (+0.2)	33.1 (-0.1)
Inside-Context: P-N-SRC	19.0 (+0.1)	19.7 (-0.1)	33.0 (-0.2)

Table 6: s-BLEU scores of the MTL and Inside-Context models are tested by giving the same source sentences as context and input. The change of s-BLEU scores over the models tested with random context is shown in ($\pm x$).

5.5 Results of MTL and Multi-Encoder models without Context

We conduct experiments on MTL and Inside-Context models by using the same source sentence as the context. Since the proposed MTL models fail when tested with random context (cf. Section 5.4), we observe how the MTL and Multi-Encoder models are performing when the same source sentence is given as context. This setting presents a scenario where the context is not random but also not the type of context with which the models are trained. We conduct experiments for *P@2-SRC* and *P-N-SRC* context settings only as the *P@2-TGT* context setting requires the current target sentence, which is unavailable during testing. We observe that MTL models can perform well compared to the random context setting, which shows that the MTL models are sensitive to the choice of context. The performance of Inside-Context models is almost the same as those tested with random context. This shows that the Inside-Context model is agnostic to the choice of the context. Table 6 shows the s-BLEU scores of the MTL and Inside-Context models.

Model	News	TED	Europarl
Vanilla-Sent	40.17	31.22	37.22
Concat-Context: P@2-SRC	39.34	30.01	36.42
Concat-Context: P-N-SRC	39.99	29.57	36.78
Concat-Context: P@2-TGT	38.50	28.82	37.27
MTL: P@2-SRC	40.69	31.44	35.96
MTL: P-N-SRC	40.50	31.24	36.94
MTL: P@2-TGT	40.99	31.90	33.91

Table 7: Accuracy of Pronoun Translation (APT) scores. The best results are shown in bold.

5.6 Pronoun Translation Accuracy

We also evaluate our proposed models’ performance on pronoun translation accuracy. We

calculate the pronoun translation accuracy with APT (accuracy of pronoun translation) (Miculicich Werlen and Popescu-Belis, 2017) metric⁶. This metric requires a list of pronouns from the source language (German) with a list of pronouns from the target language (English) as an optional argument. We use spaCy⁷ to tag both source and target sentences from the test set and extract pronouns. Table 7 shows the APT scores of *Vanilla-Sent*, *Concat-Context*, and *MTL DocNMT* models. The APT scores correlate with the s-BLEU and d-BLEU scores, achieving the highest APT score of 40.99 in *MTL: P@2-TGT* setting with an improvement of +0.82 over *Vanilla-Sent* and +1.0 over *Concat-Context (P-N-SRC)* models on News corpus. Similarly, the *MTL: P@2-TGT* model achieves the highest APT score of 31.90 with an improvement of +0.68 and +1.89 over *Vanilla-Sent* and *Concat-Context (P@2-SRC)* on TED. For the Europarl corpus, *Concat-Context (P@2-TGT)* achieved the highest APT score of 37.27 with an improvement of +0.05 and +0.33 over *Vanilla-Sent* and *MTL (P-N-SRC)* models respectively.

6 Conclusion

This work explored the MTL approach for document-level NMT (DocNMT). Our proposed MTL approach is based on cascade MTL architecture, where the model consists of one encoder (for context encoding) and two decoders (for the representation of the current source and target sentences). Reconstruction of the source sentence given the context is considered the auxiliary task, along with the translation of the current source sentence as the main task. We conducted experiments for German–English for News-commentary v14, IWSLT’17 TED, and Europarl v7 corpora with three different types of contexts *viz.* two previous sources, two previous targets, and previous-next source sentences with respect to the current input source sentence.

Our proposed MTL approaches outperform the sentence-level baseline and concatenated-context models in low-resource (for News and TED corpora) settings. However, all models perform well in the high resource setting (Europarl corpus), with proposed MTL models slightly underperforming the rest. Our MTL models are more sensitive to the choice of context than the multi-encoder mod-

els when tested with random context. We observe that the context encoder cannot encode context sufficiently and performs poorly reconstruction tasks. Finally, we reported APT (accuracy of pronoun translation) scores, and the proposed MTL models outperformed the sentence-level baseline and concatenated-context models. Our empirical analysis concludes that our approach is more sensitive to the choice of context and improves the overall translation performance in low-resource context-aware settings. We plan to explore other tasks, such as gap sentence generation (GSG) (Zhang et al., 2020a) as an auxiliary task for better context encoding, different training curricula to prioritize one objective over the other during the training, and dynamic context selection.

7 Limitations

Our study poses two main limitations. First, our primary motivation is to understand the effect of context and if the context encoding can be modelled as an auxiliary task but not to propose a model to achieve state-of-the-art results. We have followed the findings of Li et al. (Li et al., 2020) and used one of their approach to understanding the effect of context. Our observations are also in line with their findings.

Second, even though our proposed MTL approach can outperform the models in other settings, the auxiliary task (reconstruction) is not very effective as it improves the BLEU scores in the range of [0.1-1.8] over the Multi-encoder models. We hypothesize that, in the loss function, we are giving equal weights to both the objectives (0.5 for both reconstruction and translation objectives), which might lead to significantly less improvement in overall translation quality. We plan to explore different training curricula to adjust the weight of the objectives dynamically during the training.

Acknowledgements

We gratefully acknowledge the support from the “NLTM: VIDYAAPATI” project, sponsored by Electronics and IT, Ministry of Electronics and Information Technology (MeiTY), Government of India. Santanu Pal acknowledges the support from Wipro AI. We also thank the anonymous reviewers for their insightful comments.

⁶<https://github.com/idiap/APT>

⁷<https://spacy.io/models>

References

- Agrawal, Ruchit Rajeshkumar, Marco Turchi, and Matteo Negri. 2018. Contextual handling in neural machine translation: Look behind, ahead and on both sides. In *21st Annual Conference of the European Association for Machine Translation*, pages 11–20.
- Anastasopoulos, Antonios and David Chiang. 2018. Tied multitask learning for neural speech translation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 82–91, New Orleans, Louisiana, June. Association for Computational Linguistics.
- Bao, Guangsheng, Yue Zhang, Zhiyang Teng, Boxing Chen, and Weihua Luo. 2021. G-transformer for document-level machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3442–3455, Online, August. Association for Computational Linguistics.
- Barrault, Loïc, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. Findings of the 2019 conference on machine translation (WMT19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Florence, Italy, August. Association for Computational Linguistics.
- Bawden, Rachel, Rico Sennrich, Alexandra Birch, and Barry Haddow. 2018. Evaluating discourse phenomena in neural machine translation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1304–1313, New Orleans, Louisiana, June. Association for Computational Linguistics.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Donato, Domenic, Lei Yu, and Chris Dyer. 2021. Diverse pretrained context encodings improve document translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1299–1311, Online, August. Association for Computational Linguistics.
- Dong, Daxiang, Hua Wu, Wei He, Dianhai Yu, and Haifeng Wang. 2015. Multi-task learning for multiple language translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1723–1732, Beijing, China, July. Association for Computational Linguistics.
- Huo, Jingjing, Christian Herold, Yingbo Gao, Leonard Dahlmann, Shahram Khadivi, and Hermann Ney. 2020. Diving deep into context-aware neural machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 604–616, Online, November. Association for Computational Linguistics.
- Junczys-Dowmunt, Marcin. 2019. Microsoft translator at WMT 2019: Towards large-scale document-level neural machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 225–233, Florence, Italy, August. Association for Computational Linguistics.
- Kang, Xiaomian, Yang Zhao, Jiajun Zhang, and Chengqing Zong. 2020. Dynamic context selection for document-level neural machine translation via reinforcement learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2242–2254, Online, November. Association for Computational Linguistics.
- Kim, Yunsu, Duc Thanh Tran, and Hermann Ney. 2019. When and why is document-level context useful in neural machine translation? In *Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019)*, pages 24–34, Hong Kong, China, November. Association for Computational Linguistics.
- Kingma, Diederik P. and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In Bengio, Yoshua and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Koehn, Philipp. 2004. Statistical significance tests for machine translation evaluation. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain, July. Association for Computational Linguistics.
- Kudo, Taku and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium, November. Association for Computational Linguistics.

- Kudo, Taku. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia, July. Association for Computational Linguistics.
- Lei, Yikun, Yuqi Ren, and Deyi Xiong. 2022. CoDoNMT: Modeling cohesion devices for document-level neural machine translation. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 5205–5216, Gyeongju, Republic of Korea, October. International Committee on Computational Linguistics.
- Li, Bei, Hui Liu, Ziyang Wang, Yufan Jiang, Tong Xiao, Jingbo Zhu, Tongran Liu, and Changliang Li. 2020. Does multi-encoder help? a case study on context-aware neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3512–3518, Online, July. Association for Computational Linguistics.
- Liang, Yunlong, Fandong Meng, Jinan Xu, Yufeng Chen, and Jie Zhou. 2022. Scheduled multi-task learning for neural chat translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4375–4388, Dublin, Ireland, May. Association for Computational Linguistics.
- Liu, Yinhan, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.
- Luong, Minh-Thang, Quoc V Le, Ilya Sutskever, Oriol Vinyals, and Lukasz Kaiser. 2015. Multi-task sequence to sequence learning. *arXiv preprint arXiv:1511.06114*.
- Ma, Shuming, Dongdong Zhang, and Ming Zhou. 2020. A simple and effective unified encoder for document-level machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3505–3511, Online, July. Association for Computational Linguistics.
- Maruf, Sameen and Gholamreza Haffari. 2018. Document context neural machine translation with memory networks. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1275–1284, Melbourne, Australia, July. Association for Computational Linguistics.
- Maruf, Sameen, André F. T. Martins, and Gholamreza Haffari. 2019. Selective attention for context-aware neural machine translation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3092–3102, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Miculicich, Lesly, Dhananjay Ram, Nikolaos Pappas, and James Henderson. 2018. Document-level neural machine translation with hierarchical attention networks. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2947–2954, Brussels, Belgium, October–November. Association for Computational Linguistics.
- Miculicich Werlen, Lesly and Andrei Popescu-Belis. 2017. Validation of an automatic metric for the accuracy of pronoun translation (APT). In *Proceedings of the Third Workshop on Discourse in Machine Translation*, pages 17–25, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Popel, Martin and Ondřej Bojar. 2018. Training tips for the transformer model. *arXiv preprint arXiv:1804.00247*.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October. Association for Computational Linguistics.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August. Association for Computational Linguistics.
- Tiedemann, Jörg and Yves Scherrer. 2017. Neural machine translation with extended context. In *Proceedings of the Third Workshop on Discourse in Machine Translation*, pages 82–92, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Voita, Elena, Pavel Serdyukov, Rico Sennrich, and Ivan Titov. 2018. Context-aware neural machine translation learns anaphora resolution. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2947–2954, Brussels, Belgium, October–November. Association for Computational Linguistics.

- of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1264–1274, Melbourne, Australia, July. Association for Computational Linguistics.
- Voita, Elena, Rico Sennrich, and Ivan Titov. 2019. When a good translation is wrong in context: Context-aware machine translation improves on deixis, ellipsis, and lexical cohesion. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1198–1212, Florence, Italy, July. Association for Computational Linguistics.
- Wang, Yiren, ChengXiang Zhai, and Hany Hassan. 2020. Multi-task learning for multilingual neural machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1022–1034, Online, November. Association for Computational Linguistics.
- Wang, Tao, Chengqi Zhao, Mingxuan Wang, Lei Li, and Deyi Xiong. 2021. Autocorrect in the process of translation — multi-task learning improves dialogue machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Papers*, pages 105–112, Online, June. Association for Computational Linguistics.
- Wu, Yonghui, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144.
- Yang, Jiacheng, Mingxuan Wang, Hao Zhou, Chengqi Zhao, Weinan Zhang, Yong Yu, and Lei Li. 2020. Towards making the most of bert in neural machine translation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 9378–9385.
- Zaremoondi, Poorya and Gholamreza Haffari. 2018. Neural machine translation for bilingually scarce scenarios: a deep multi-task learning approach. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1356–1365, New Orleans, Louisiana, June. Association for Computational Linguistics.
- Zaremoondi, Poorya, Wray Buntine, and Gholamreza Haffari. 2018. Adaptive knowledge sharing in multi-task learning: Improving low-resource neural machine translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 656–661, Melbourne, Australia, July. Association for Computational Linguistics.
- Zhang, Jiacheng, Huanbo Luan, Maosong Sun, Feifei Zhai, Jingfang Xu, Min Zhang, and Yang Liu. 2018. Improving the transformer translation model with document-level context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 533–542, Brussels, Belgium, October-November. Association for Computational Linguistics.
- Zhang, Jingqing, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020a. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In *International Conference on Machine Learning*, pages 11328–11339. PMLR.
- Zhang, Pei, Boxing Chen, Niyu Ge, and Kai Fan. 2020b. Long-short term masking transformer: A simple but effective baseline for document-level neural machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1081–1087, Online, November. Association for Computational Linguistics.
- Zhou, Shuyan, Xiangkai Zeng, Yingqi Zhou, Antonios Anastasopoulos, and Graham Neubig. 2019. Improving robustness of neural machine translation with multi-task learning. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 565–571, Florence, Italy, August. Association for Computational Linguistics.

A Appendix

A.1 Preliminary Investigation on Auxiliary Objectives

The joint probability in Equation 2 ($p(x, c_x)$) can be calculated in two ways such as:

$$p(x, c_x) = p(x|c_x) \times p(c_x) \quad (6)$$

$$p(x, c_x) = p(c_x|x) \times p(x) \quad (7)$$

Since the joint probability can be computed in two different ways, we conduct an initial study to select the optimal auxiliary objective that improves the overall translation performance of the model. Specifically, we consider $p(x|c_x)$ as one auxiliary task where source (x) is autoregressively reconstructed (denoted as **Re-Src**) from the encoded context (c_x) and $p(c_x|x)$ as the other auxiliary task where context (x) is autoregressively

reconstructed (denoted as **Re-Cntx**) from the encoded source (c_x). We conducted experiments to verify which auxiliary task is performing better.

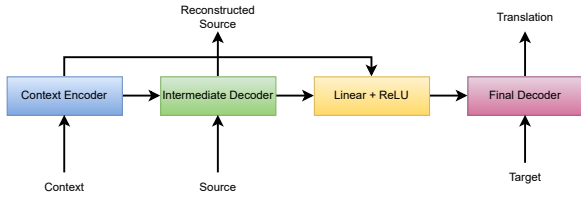


Figure 3: The overview of modified MTL architecture with residual connection. The input to the model is a triplet. The triplet consist of (*Context*, *Source*, *Target*) in **Re-Src** setting and (*Source*, *Context*, *Target*) in **Re-Cntx** setting. Here, *Source*: Current source sentence, *Context*: Context for the current source sentence, and *Target*: Translation of current source sentence. None of the layers are shared.

The experimental setup and model architecture are slightly different for this comparison study than those used in the main experiments.⁸ The Context Encoder and Intermediate Decoder output are combined with a linear layer with ReLU activation. The main experimental setup does not use this linear layer + ReLU combination. We hypothesize that adding this layer might make the model agnostic to the choice of context. We test this by training the model with random context (cf. Section A.1.1). Specifically, we use two context settings *viz.* *P@2-SRC* and *P-N-SRC* settings (cf. 3.3). We use a fixed learning rate of 10^{-5} instead of the warmup schedule. The output from this layer is given as input to the Final Decoder.

Model		Vanilla-Sent	MTL: P@2-SRC	MTL: P-N-SRC
News	Re-Src		20.6	20.9
	Re-Cntx	16.5	16.7 (-3.9)	17.9 (-3.0)
TED	Re-Src		21.6	22.0
	Re-Cntx	12.1	18.0 (-3.6)	17.8 (-4.2)
Europarl	Re-Src		35.1	35.8
	Re-Cntx	35.0	33.2 (-1.9)	33.6 (-2.2)

Table 8: Comparison of s-BLEU scores of Baseline and proposed MTL DocNMT models trained with different source contexts for German to English direction. Differences in the scores over **Re-Src** are shown inside the parentheses.

We use a mini-batch of 18 sentences to train all the models. We create two separate subword vocabularies for each training corpus. The created subword vocabulary is 40k in both German and English. We use the unigram language model

⁸We modified the experimental setup and model architecture during our main experiments. In this preliminary investigation, the capacity of models with independent subword vocabularies is slightly larger. Due to this, the s-BLEU scores are slightly better than the main results.

(Kudo, 2018) to create subword vocabularies with SentencePiece (Kudo and Richardson, 2018), and the maximum sequence length is set to 160 tokens. During inference, we perform greedy decoding. The rest of the experimental setup is the same as the one used in the main experiments.

Model	Random-Train		Random-Infer	
	Re-Src	Re-Cntx	Re-Src	Re-Cntx
MTL: P@2-SRC	20.9	16.6	20.6	16.8
MTL: P-N-SRC	20.9	16.4	20.8	17.8

Table 9: s-BLEU scores of *Random-Train* and *Random-Infer* experiments on News-commentary corpus.

A.1.1 Effect of Random Context

We also conduct experiments to study how the random context affects the MTL models. Specifically, we evaluate the MTL models in two settings. The model is trained on the random context in the *Random-Train* setting by concatenating two randomly sampled sentences from the train set and testing with *P@2-SRC* and *P-N-SRC* context settings. In *Random-Infer* setting, the model is trained on *P@2-SRC* and *P-N-SRC* context settings and tested with random context. We train the models on the news-commentary corpus. Table 9 shows the s-BLEU scores of the MTL models trained and tested in the random context setting. Based on the results, we conclude that the model trained with random context improves the robustness of the model. This observation aligns with the findings of Li et al. (2020), but they conducted experiments in the non-MTL setting with multiple encoders. As the model largely ignores the choice of the context, we remove this linear + ReLU combination and feed the output of the Intermediate Decoder to the Final Decoder. We hypothesize that this forces the model to consider the context while generating the target sentence.

Aligning Neural Machine Translation Models: Human Feedback in Training and Inference

Miguel Moura Ramos^{1,2} Patrick Fernandes^{1,2,3} António Farinhas^{1,2}
André F. T. Martins^{1,2,4}

¹Instituto Superior Técnico, Universidade de Lisboa (ELLIS Unit Lisbon)

²Instituto de Telecomunicações ³Carnegie Mellon University ⁴Unbabel

miguel.moura.ramos@tecnico.ulisboa.pt

Abstract

Reinforcement learning from human feedback (RLHF) is a recent technique to improve the quality of the text generated by a language model, making it closer to what humans would generate. A core ingredient in RLHF’s success in aligning and improving large language models (LLMs) is its *reward model*, trained using human feedback on model outputs. In machine translation (MT), where metrics trained from human annotations can readily be used as reward models, recent methods using *minimum Bayes risk* decoding and reranking have succeeded in improving the final quality of translation. In this study, we comprehensively explore and compare techniques for integrating quality metrics as reward models into the MT pipeline. This includes using the reward model for data filtering, during the training phase through RL, and at inference time by employing reranking techniques, and we assess the effects of combining these in a unified approach. Our experimental results, conducted across multiple translation tasks, underscore the crucial role of effective data filtering, based on estimated quality, in harnessing the full potential of RL in enhancing MT quality. Furthermore, our findings demonstrate the effectiveness of combining RL training with reranking techniques, showcasing substantial improvements in translation quality.

1 Introduction

Neural machine translation (NMT) models (Bahdanau et al., 2015; Vaswani et al., 2017) are typically trained with *maximum likelihood estimation*

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

(MLE), maximizing the log-probability of the next word in a translation given the previous words and the source sentence. While this approach has been effective at training high-quality MT systems, the difference between the training and inference objective can lead to *exposure bias* (Bengio et al., 2015; Ranzato et al., 2016; Wiseman and Rush, 2016), which hinders the model’s ability to recover from early mistakes. Furthermore, the suitability of model likelihood as a proxy for generation quality has been questioned in machine translation (Koehn and Knowles, 2017; Ott et al., 2018) and beyond (Perez et al., 2022). These challenges sparked interest in alternative training and decoding paradigms for MT, such as *reinforcement learning* (RL; Kreutzer et al. (2018)) or *minimum Bayes risk* decoding (MBR; Eikema and Aziz (2022)).

More recently, the widespread success of *reinforcement learning from human feedback* (Stienon et al., 2022) has highlighted the importance of a good reward model that approximates well to human preferences for the task at hand. While, in general, this requires training a reward model from scratch for the specific problem, in the case of machine translation (MT), the evaluation community has achieved significant progress in developing automatic quality estimation and evaluation metrics *learned* from human quality annotations (e.g. COMET-QE (Rei et al., 2020), COMET (Rei et al., 2022a), BLEURT (Sellam et al., 2020), which can be repurposed as reward models. As a consequence, recent research integrating these metrics into the training (Gulcehre et al., 2023) or decoding (Fernandes et al., 2022) procedures has had considerable success in improving the quality of translations. However, none of the previous work has systematically compared the effect of integrating metrics at different stages of the MT pipeline or has attempted to combine these techniques in a unified approach.

In this work, we perform a comprehensive study

on the integration of MT quality metrics into the MT pipeline as reward models. As illustrated in Figure 1, we assess their use at different stages: as a means for data filtering, during the training process through RL, and at inference time by way of reranking techniques. Furthermore, we explore the results of combining these methods.

We attempt to answer the following research questions:

- *Can data filtering based on estimated quality help minimize RL training instability?*
- *Which metrics are more suitable as reward models in RL training? Are reference-free metrics competitive with reference-based ones?*
- *How does the quality of translations achieved through RL training compare with those produced through reranking approaches? Can these two approaches be effectively combined to further enhance translation quality?*

Our main contributions arise from the research questions mentioned above:

- Inspired by Bane and Zaretskaya (2021) where they use cross-lingual encoders to score translation representations in an aligned multilingual vector space, we propose an alternative data filtering method that uses COMET-QE (Rei et al., 2020), a more robust model, to curate a high-quality dataset that empirically helps to **minimize RL training instability**.
- We show that neural metrics such as COMET(-QE) (Rei et al., 2022a; Rei et al., 2020) are more suitable than BLEU (Papineni et al., 2002) for RL training. Contrary to what happens with MBR decoding, RL training results in improved scores across all types of metrics, not only neural ones. In particular, using a reward model based on QE works surprisingly well, possibly paving the way for unsupervised training of NMT systems.
- Experiments in EN→DE and EN→FR show that both RL training and reranking techniques enhance translation quality, with RL training often outperforming reranking methods. Furthermore, combining RL and MBR decoding results in more consistent improvements across various evaluation metrics.

- We quantify and discuss the trade-offs in running time at both training and inference, clarifying the efficiency and suitability of each approach.

2 Background

2.1 Neural Machine Translation

An NMT model has learnable parameters, θ , to estimate the probability distribution, $p_\theta(y|x)$ over a set of hypotheses \mathcal{Y} , conditioned on a source sentence x . MLE is the training principle of estimating θ , given parallel data, formalized as

$$\mathcal{L}(\theta, y_{1:L}) = -\frac{1}{L} \sum_{t=1}^L \log p_\theta(y_t | y_0, \dots, y_{t-1}). \quad (1)$$

NMT systems typically employ *maximum a posteriori* (MAP) decoding to generate translations,

$$\hat{y}_{\text{MAP}} = \arg \max_{y \in \mathcal{Y}} \log p_\theta(y|x), \quad (2)$$

where algorithms such as greedy decoding or beam search (Reddy, 1977) approximate the most *probable* translation given the source. An alternative approach is to sample translations according to $p_\theta(y|x)$, using techniques such as top- k or nucleus sampling (Fan et al., 2018; Holtzman et al., 2020).

In §3.3 of this paper, we also consider two distinct reranking approaches (Fernandes et al., 2022), namely N -best reranking and MBR decoding. While N -best reranking selects the candidate translation that maximizes a given (reference-free) metric, MBR decoding ranks candidates using reference-based metrics, maximizing the expected utility (or minimizing the risk).

2.2 MT Evaluation

Human evaluations are the most reliable way to assess the performance of MT systems, but they are time-consuming and costly. For that reason, the standard way to evaluate MT is through automatic evaluation metrics, which can be reference-based or quality estimation (QE) metrics.

Reference-based metrics compare the generated translation to human-written reference texts. Lexical reference-based metrics, such as the widely used BLEU (Papineni et al., 2002), rely on word overlap and n-gram matching, making them ineffective for translations that have the same meaning but are substantially different from the reference. On the other hand, neural metrics, such as COMET (Rei et al., 2022a), are a recent alternative that relies on neural networks trained on human-annotated

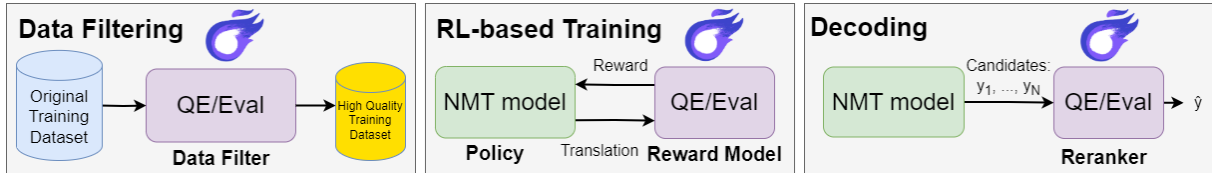


Figure 1: Preference models can have multifaceted roles within the MT pipeline. They can serve as effective data filters, refining datasets by incorporating user preferences. They can also assume a pivotal role in classic RL training by providing rewards to optimize the MT model performance. Finally, they can act as rerankers during the decoding phase, selecting the final translation by maximizing their scores derived from user preferences.

data and that leverages contextual embeddings to address semantic similarity.

QE assesses translation quality without human references, being particularly useful in dynamic, data-intensive environments, where references are costly and time-consuming to obtain. This paper focuses on sentence-level QE as a reward model, providing a single quality assessment for each translation. COMET-QE (Rei et al., 2020) is a state-of-the-art reference-free quality estimation metric derived from COMET used to evaluate MT performance.

Neural reference-based and QE metrics are valuable preference models because they offer a more accurate and contextually-aware measure of translation quality, aligning better with human preferences and judgments (Freitag et al., 2022b).

2.3 Reinforcement Learning Training in NMT

In MT, approaches based on reinforcement learning (RL; Sutton and Barto (2018)) cast the problem as a Markov decision process (MDP; Puterman (1990)), where a source sentence $x = (x_1, \dots, x_n)$ is translated into a target sentence $y = (y_1, \dots, y_m)$. Under this perspective, the NMT system can be viewed as the agent with a conditional probability distribution based on its parameters, $p_\theta(y_t|x, y_{<t})$. The states of the MDP are defined by the target sentence that has already been decoded, $s_t = (y_1, \dots, y_{t-1})$, and the action corresponds to the selection of the next word, y_{t+1} . Based on the states and actions, all transitions are deterministic and the reward function, R , is provided by the MT evaluation model which returns a quality score for the generated translation \hat{y} . The main purpose of using RL in NMT is to provide learning signals that go beyond a single reference translation, by providing reward signals for arbitrary translations. MLE provides less robust learning signals that are more susceptible to the shortcomings of noisy references. However, it

is essential to note that if the reward model used relies on reference-based metrics, some vulnerability to noisy references may still persist. Accordingly, the goal of RL training is to maximize the expected reward, $L_{rl}(\theta) = \mathbb{E}_{p_\theta(\hat{y}|x)}[R(\hat{y})]$. Commonly used RL training procedures include REINFORCE (Williams, 1992), minimum risk training (Och, 2003; Shen et al., 2016), and proximal policy optimization (PPO; Schulman et al. (2017)).

3 Aligning MT with Reward Models

3.1 Data Filtering

The success of fine-tuning NMT models with MLE is highly dependent on the quantity and quality of the training dataset (Wang et al., 2018; Koehn and Knowles, 2017; Khayrallah and Koehn, 2018). This is because accurate references are crucial for computing meaningful learning signals that correctly guide the NMT model towards improved translations (Kong et al., 2018). Despite its recent successes, RL-based training can be unstable, so using only high-quality data could help mitigate this instability. This can be addressed via **data filtering**, by seeking a good balance between the aggressiveness of filtering and the resulting dataset size: if the original dataset is already small, too much filtering can be detrimental to the performance of NMT systems (Zoph et al., 2016; Jiao et al., 2020). Furthermore, when looking at the RL scenario, having a sufficiently large training dataset can help guarantee that the NMT model explores a wide range of scenarios for policy improvement.

We apply our data filtering method on the considerably large and noisy WMT datasets (Bojar et al., 2015; Bojar et al., 2016) since they have been reported to have less relevant and uncorrelated sentences that can lead to sub-optimal results when used during training (Koehn et al., 2020; Malli and Tambouratzis, 2022). We do not perform data filtering to the IWSLT2017 (Cettolo et al., 2012;

Cettolo et al., 2017) dataset due to concerns about its limited amount of available data. Further dataset filtering could potentially result in a too-small training dataset, which is not desirable for training MT systems.

As illustrated in Figure 1, to perform the training dataset filtering, we use a filter that reranks the sentence pairs according to quality scores that indicate the correlation and relevance of each sentence and its given reference. This approach allows us to filter out low-quality sentence pairs, thereby improving the overall quality of the data. In our approach, we use a robust preference model called COMET-QE (Rei et al., 2020) as the data filter, which combines the use of encoders and a regression model trained on human-annotated data to estimate the quality score of each sentence pair. This reference-less model is expected to be more accurate in quality score estimation and have a superior alignment with human judgments than just resorting to the currently used cross-lingual encoders which only take into account vector-space mapping similarity (Bane and Zaretskaya, 2021). Furthermore, COMET-QE seems particularly suitable as our preference model during data filtering, as it is a multilingual reference-free neural-based metric trained on human annotations of translation quality, and therefore can be used to filter by thresholding on predicted quality or on the number of sentences in the training set. After scoring all sentence pairs, we select the threshold based on the number of high-quality sentence pairs to use as the filtered dataset for RL training. For that, we apply different thresholds and sizes to the reranked sentences. We, then, MLE fine-tune our baseline on these subsets and select the subset that gives the overall best-performing model on the dev. set. These best-performing models serve as baselines for our RL-based training and reranking methods during decoding.

In conclusion, it is worth noting that our data filtering method is, as shown in Figure 1, one of three methods we cover for employing a preference model in the MT pipeline. This filtering method can significantly increase the performance of MT systems by introducing feedback in an earlier stage of the pipeline.

3.2 Training Phase

The use of RL-based training has the potential to bridge the gap between MLE training objectives, MT evaluation metrics and human-like translations.

However, it faces challenges of instability and inefficiency, especially in gradient variance and reward computation. As illustrated in Figure 1, the RL training process is composed of an NMT model that generates translations that are evaluated by the reward model through rewards that represent the quality of the translation. This reward is used by the policy gradient algorithm to update the NMT model’s **policy**. To address the problem of gradient variance, we employ PPO (Schulman et al., 2017) as our policy gradient algorithm since it is a stable and efficient algorithm that updates the policy parameters in a controlled way with a predetermined proximity bound, avoiding sudden changes that might destabilize the learning.

Reward computation is the most crucial part of this entire process as it guides the NMT model during training. Previous work on RL-based NMT systems predominantly used BLEU as the reward function. However, BLEU has several limitations, as discussed in §2.2. To address these shortcomings, we leverage robust preference models during RL training, such as the reference-based COMET (Rei et al., 2022a) and the reference-free COMET-QE (Rei et al., 2020), as highlighted in Figure 1. Since learning these models is a complex task, we incorporate these pre-trained preference models, which have already been shown to correlate well with human judgments (Freitag et al., 2022b; Rei et al., 2022a; Rei et al., 2020), to ensure that RL systems can better capture the nuanced preferences of the user by receiving human-like feedback as rewards. These models assign numerical quality scores to each translation hypothesis based on their desirability, making them similar to utility functions. Our study aims to demonstrate that training with RL can generate higher-quality NMT models using neural metrics and investigate the competitiveness of COMET-QE as a reward model.

Another crucial decision was related to the exploitation vs. exploration problem of RL in the context of MT (Wu et al., 2018). The beam search algorithm generates more accurate translations by exploiting the probability distribution/policy of the NMT model, while sampling aims to explore more diverse candidates. During generation, we observed that sampling techniques generally led to candidates of lower quality when compared to beam search, according to the preference models used. Therefore, all RL-based models used beam search during their training and inference.

3.3 Decoding Phase

Reranking methods (Ng et al., 2019; Bhattacharyya et al., 2021; Fernandes et al., 2022; Eikema and Aziz, 2022) are an alternative to MAP-based decoding that relies on reranking techniques and presupposes access to N candidate translations for each source sentence, generated by the NMT system through methods like beam search or sampling. The generated candidates are reranked according to their quality given an already determined metric/reward model.

We employ two reranking methods to select a final translation: N -best reranking (Ng et al., 2019; Bhattacharyya et al., 2021) and *minimum Bayes risk* decoding (MBR; Eikema and Aziz (2022)).

N -best reranking (3) employs a reference-free metric, M_{QE} , to reorder a set of N candidate translations, denoted as $\bar{\mathcal{Y}}$, and selects the candidate with the highest estimated quality score as the final translation, \hat{y}_{RR} ,

$$\hat{y}_{RR} = \arg \max_{y \in \bar{\mathcal{Y}}} M_{QE}(y). \quad (3)$$

Considering the previous equation, and assuming $C_{M_{QE}}$ as the computational cost of evaluating a candidate translation with QE metric, M_{QE} , we obtain the final computational cost of finding the best translation from N candidate translations as $O(N \times C_{M_{QE}})$.

MBR decoding, in contrast, relies on a reference-based metric and chooses the candidate that has the highest quality when compared to other possible translations (in expectation). We define $u(y^*, y)$ as the utility function, quantifying the similarity between a hypothesis $y \in \mathcal{Y}$ and a reference $y^* \in \bar{\mathcal{Y}}$. In our context, the utility function is represented by either BLEU or COMET. Therefore, MBR decoding can be mathematically expressed as

$$\hat{y}_{MBR} = \arg \max_{y \in \bar{\mathcal{Y}}} \underbrace{\mathbb{E}_{Y \sim p_\theta(y|x)}[u(Y, y)]}_{\approx \frac{1}{N} \sum_{j=1}^N u(y^{(j)}, y)}, \quad (4)$$

where in Eq. 4 the expectation is approximated as a Monte Carlo sum using model samples $y^{(1)}, \dots, y^{(N)} \sim p_\theta(y|x)$. These samples may be obtained through biased sampling (e.g., nucleus-p or top-k) or beam search. Knowing that the utility function is a reference-based metric M_{REF} with computational cost, $C_{M_{REF}}$, and that to find the best translation we need to do pairwise comparisons between hypotheses, we obtain the final computational cost as $O(N^2 \times C_{M_{REF}})$. These reranking

methods become particularly effective when N is not excessively large, making the process computationally more manageable.

Preference models capture the preferences of human evaluators and can be used during the decoding stage to influence MT systems, as shown in Figure 1. By doing this, the MT system will prioritize translations that are more aligned with human judgments, therefore reducing the chances of generating severely incorrect translations. We believe that incorporating preference models during the decoding stage can lead to even better translation quality, even if the underlying model has already been RL-trained using the same or a different preference model. The benefits we expect to see include improved fluency, adequacy, and consistency compared to the respective baselines since our preference models have been trained on annotations that aim to optimize these linguistic aspects.

4 Experiments

4.1 Setup

During the training phase, we investigate the advantages of RL training (with and without data filtering §3.1) for enhancing the performance of NMT systems. We employ a T5 model¹, pre-trained on the C4 dataset (Raffel et al., 2019). First, we fine-tune the models using MLE training with Adam (Kingma and Ba, 2017) as the optimization algorithm, learning rate decay starting from 5×10^{-6} and early stopping. For RL training², we use PPO with learning rate set as 2×10^{-5} , γ set as 0.99, trajectory limit set as 10,000, beam search size set as 5 and mini-batch updates were conducted using stochastic gradient descent with a batch size of 32, gathered over 4 PPO epochs. In the inference phase, our emphasis shifts towards reranking techniques and their impact on the performance of NMT systems. As for the candidate generation method used, early experiments, omitted for relevancy, show that the best configuration is to generate 100 candidates per source sentence and then use sampling with $p = 0.6$ and $k = 300$ to select the best translation. Consequently, the evaluation encompasses all the baseline and RL-trained models, both with and without N -best reranking and MBR decoding. These evaluations are conducted across the follow-

¹We leverage the T5-Large model available in Huggingface’s *Transformers* framework (Wolf et al., 2020).

²Our RL implementation relies on the Transformer Reinforcement Learning X framework (Castricato et al., 2023, trIX).

ing datasets:

- The small IWSLT2017 datasets (Cettolo et al., 2012; Cettolo et al., 2017) for English to German (EN → DE) and English to French (EN → FR), featuring 215k and 242k training examples, respectively.
- The large WMT16 dataset (Bojar et al., 2016) for English to German (EN → DE) with 4.5M training examples.
- The large WMT15 dataset (Bojar et al., 2015) for English to French (EN → FR) with over 40M training samples.

We assess the performance of each NMT system using well-established evaluation metrics, which include BLEU, chrF (Popović, 2015), METEOR (Banerjee and Lavie, 2005), COMET, COMET-QE, and BLEURT. Additionally, for certain experiments executed on a single NVIDIA RTX A6000 GPU, we provide wall clock time measurements to offer insights into computational efficiency.

4.2 Finding the Optimal Quality Subset Size

In this section, we discuss our approach to quality-aware data filtering as a stabilizing strategy (§3.1), for the WMT datasets. Figure 2a summarizes our findings for the WMT16 EN→DE dataset (Bojar et al., 2016) on the influence of a high-quality subset on translation performance as we vary the subset size, based on various evaluation metrics and COMET-QE sentence filtering. Across all metrics, a consistent trend emerges: after reaching training sizes of 500 000, there is a notable decline in performance. Particularly, this decline is less prominent for lexical metrics, possibly due to their inherent limitations (Freitag et al., 2022b). A similar analysis for WMT15 EN→FR that can be found in Figure 2b results in an optimal training size of 300 000 examples.

While the data filtering process has led to remarkable improvements in performance, it is important to note that the effectiveness of this process is dependent on the selected reranking metric. Using metrics that are not closely aligned with human judgments can result in poorly correlated and misaligned sentences, which can make the training process more unstable. Therefore, it is recommended to use robust QE models, such as COMET-QE. The more recent COMETKIWI (Rei et al., 2022b) model may offer even greater performance improvements.

4.3 Impact of Quality-aware Data Filtering

After obtaining the best configuration for our data filtering process, we experiment with the use of the curated high-quality training subset from COMET-QE and assess its impact on the MLE and RL training performance. We compare our filtering method with no filtering by using the original full training dataset, random filtering and cross-lingual embedding similarity filtering using MUSE (Lample et al., 2017) and XLM-R (Conneau et al., 2019).

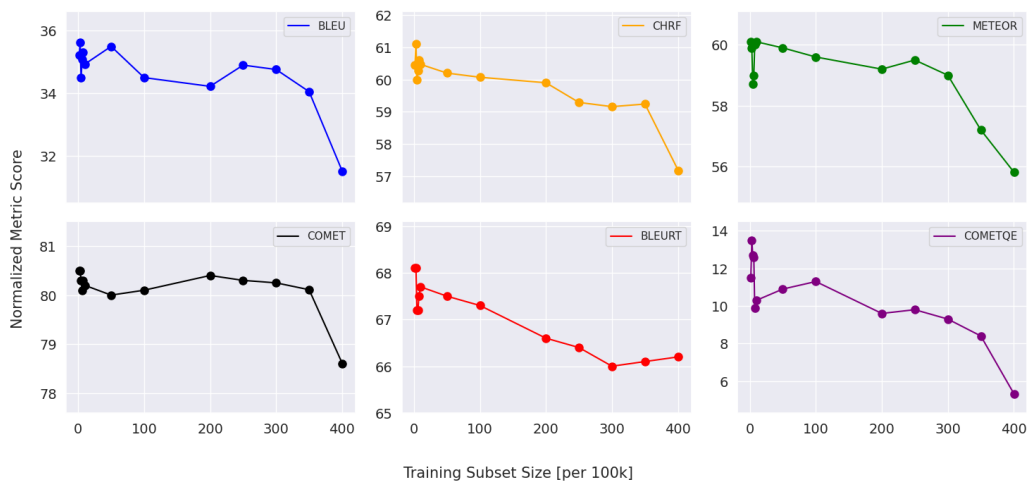
Table 1 provides a comprehensive overview of the experimental results using BLEU, COMET and COMET-QE as reward models. Both MT tasks demonstrate the same tendency when trained using MLE. COMET-QE and MUSE high-quality subsets have enough reduced noise to provide more stable training, as evidenced by the overall increase in performance across all metrics compared to the baseline training on the full original dataset. Moreover, a randomly selected subset fine-tuned with MLE performs worse or at most not significantly better than the baseline trained on the original dataset, as expected. Furthermore, in accordance with our expectations (Bane and Zaretskaya, 2021), XLM-R filtering does not improve training and is actually the worst-performing model.

Regarding RL-based training on both MT tasks, we observe that most RL-trained models outperform their MLE-trained baseline counterparts across various metrics. Notably, the best-performing models are the ones that were MLE fine-tuned and then RL-trained on the COMET-QE high-quality subset using both COMET and COMET-QE as reward models. On top of that, we can see that in some cases RL training solely does not yield significant improvements, but when combined with high-quality training subsets, it results in substantial enhancements and a competitive edge over the normal, random and XLM-R baselines. Additionally, we see impressive BLEU scores with RL training with COMET(-QE) as reward model. This finding underscores that optimizing for COMET(-QE) yields superior BLEU scores compared to direct optimization for BLEU. This phenomenon is likely attributed to COMET(-QE) providing **more effective** reward signals during training, thus highlighting the limitations of BLEU.

The excellent performance gains with COMET-QE as a data filter and also as a reward model emphasize the potential of RL-based NMT models trained with a QE reward model (which *does not* re-



(a) Impact of Data Filtering on WMT16 En→De



(b) Impact of Data Filtering on WMT15 En→FR

Figure 2: These models were fine-tuned by progressively increasing the size of the high-quality subset, obtained with COMET-QE sentence reranking and denoted in increments of 100,000.

quire a corpus with references) to outperform other RL-trained models, offering promising opportunities for unsupervised NMT training with monolingual data, especially for low-resource languages, by eliminating the need for reference translations in evaluation and reward signal generation.

In conclusion, we highlight the importance of thoughtful data selection for achieving better translation quality, showing that COMET-QE can consistently outperform the remaining filtering methods. Furthermore, the top-performing models were RL-trained with **neural** metrics, showing once again that **human-aligned** preference models can constantly outperform simpler metrics, such as BLEU.

4.4 Impact of preference-based MT alignment

Table 2 presents the performance scores of the best baseline model, across various MT tasks, focusing on the comparison between RL training, reranking

methods during inference and the potential synergies between RL training and reranking techniques in improving the translation quality of MT systems.

Our analysis reveals consistent improvements across all evaluation metrics and reward models, with RL training consistently achieving top scores, especially when using COMET-QE as the reward model.³ MBR decoding with COMET and *N*-best reranking with COMET-QE outperformed RL training in COMET and COMET-QE metrics but had difficulty improving other evaluation metrics, while RL training exhibited better generalization with slightly less consistent improvements in COMET and COMET-QE scores. This phenomenon of increased COMET and COMET-QE scores comes at the cost of worse performance according to the

³We also provide additional fine-grained quality analysis in Appendix A to better illustrate and address specific research questions.

Training Data		Lexical Metrics			Neural Metrics		
SL Data	RL Data	BLEU	ChrF	METEOR	COMET	COMET-QE	BLEURT
MLE							
Original	-	35.04	61.30	61.91	84.40	39.50	74.70
Random	-	34.43	61.00	61.36	83.90	39.10	74.30
XLM-R	-	33.24	60.35	60.20	84.80	41.80	72.60
MUSE	-	35.10	61.90	62.20	85.10	40.40	74.30
COMET-QE	-	<u>35.45</u>	<u>62.00</u>	<u>62.75</u>	<u>85.50</u>	<u>42.00</u>	<u>75.90</u>
RL w/ BLEU							
Original	Original	34.70	60.90	61.45	85.60	42.20	74.60
Random	Random	34.49	61.10	61.49	85.60	42.20	74.40
XLM-R	XLM-R	33.21	60.41	60.10	85.10	42.70	73.10
MUSE	MUSE	35.34	62.10	62.73	85.60	40.80	74.50
Original	COMET-QE	35.37	61.70	62.04	85.40	41.00	74.20
COMET-QE	COMET-QE	<u>35.55</u>	<u>62.10</u>	<u>62.77</u>	86.80	45.00	76.10
RL w/ COMET							
Original	Original	35.05	61.30	61.82	85.60	41.80	74.40
Random	Random	34.96	61.40	61.80	85.60	41.80	74.20
XLM-R	XLM-R	33.60	60.74	60.40	85.00	42.00	72.90
MUSE	MUSE	35.18	61.90	62.56	85.50	41.90	74.60
Original	COMET-QE	35.58	61.80	62.20	85.70	41.70	74.50
COMET-QE	COMET-QE	<u>35.90</u>	<u>62.20</u>	<u>63.06</u>	<u>86.70</u>	<u>44.10</u>	<u>75.70</u>
RL w/ COMET-QE							
Original	Original	34.21	60.50	61.10	85.60	42.40	74.80
Random	Random	34.88	61.30	61.69	85.50	41.80	74.10
XLM-R	XLM-R	33.57	60.73	60.40	85.10	42.20	73.20
MUSE	MUSE	35.03	61.90	62.57	85.70	41.30	74.70
Original	COMET-QE	35.48	61.70	62.10	85.70	41.70	74.50
COMET-QE	COMET-QE	<u>35.96</u>	<u>62.30</u>	<u>63.07</u>	<u>86.70</u>	<u>44.70</u>	<u>75.90</u>
WMT15 EN-FR							
Training Data		Lexical Metrics			Neural Metrics		
SL Data	RL Data	BLEU	ChrF	METEOR	COMET	COMET-QE	BLEURT
MLE							
Original	-	31.49	57.18	55.80	78.60	5.30	66.20
Random	-	31.27	57.07	60.01	80.00	12.80	65.20
XLM-R	-	25.04	48.78	48.60	77.40	12.10	57.10
MUSE	-	35.49	59.10	60.55	80.10	13.10	67.50
COMET-QE	-	<u>35.62</u>	<u>59.90</u>	<u>61.11</u>	<u>80.50</u>	<u>13.50</u>	<u>68.10</u>
RL w/ BLEU							
Original	Original	35.47	59.90	61.03	80.20	16.90	67.10
Random	Random	32.75	58.10	60.20	80.03	14.10	66.35
XLM-R	XLM-R	25.78	49.69	49.30	77.70	13.30	57.80
MUSE	MUSE	35.55	60.10	60.56	81.90	17.10	67.50
Original	COMET-QE	35.67	60.10	61.01	81.20	17.10	67.30
COMET-QE	COMET-QE	<u>36.26</u>	<u>60.40</u>	<u>61.51</u>	<u>82.10</u>	<u>17.50</u>	<u>67.70</u>
RL w/ COMET							
Original	Original	35.50	59.90	61.00	80.40	16.80	67.00
Random	Random	34.15	59.50	60.93	80.50	15.50	67.10
XLM-R	XLM-R	25.08	48.84	48.60	77.50	12.40	57.20
MUSE	MUSE	36.00	60.10	61.20	80.80	17.00	67.30
Original	COMET-QE	35.98	60.00	61.09	81.80	17.10	67.20
COMET-QE	COMET-QE	<u>36.62</u>	<u>60.60</u>	<u>61.79</u>	<u>82.20</u>	<u>17.40</u>	<u>67.60</u>
RL w/ COMET-QE							
Original	Original	35.50	60.00	61.10	82.20	17.50	68.00
Random	Random	32.10	58.30	60.50	81.00	14.40	66.70
XLM-R	XLM-R	24.67	48.38	48.10	77.60	12.60	56.80
MUSE	MUSE	35.62	60.45	59.30	82.22	17.45	67.80
Original	COMET-QE	35.90	60.10	61.22	82.27	17.53	68.02
COMET-QE	COMET-QE	<u>36.25</u>	<u>60.50</u>	<u>61.58</u>	82.40	17.70	68.10

Table 1: Automatic evaluation metrics for the MLE and RL-trained models on the WMT16 EN→DE (top) and WMT15 EN-FR (bottom) original datasets, quality subsets obtained from COMET-QE, XLM-R and MUSE and a randomly selected subset. The training data used for MLE and RL can be found in the SL and RL Data, respectively. We experimented with BLEU, COMET and COMET-QE as reward models for the RL training. The best overall values are **bolded** and the best for each specific group are underlined.

other MT evaluation metrics, showing a potential of **overfitting** effect for these reranking techniques that occur across all datasets. These findings underscore the potential of neural metrics as reward signals in training and inference, as discussed in Deutsch et al. (2022) and Freitag et al. (2022b). While combining RL training and MBR decoding occasionally led to top performance, it did not consistently outperform other strategies, making it

a method that distributes gains across all evaluation metrics without exceptional generalization as RL training but provides better overall scores than reranking methods alone.

RL training and MBR decoding in MT exhibit distinct computational efficiency profiles, as shown in Table 3. RL training is computationally demanding but typically entails a one-time, resource-intensive training process (though less resource-

MODEL	WMT16 EN→DE						WMT15 EN→FR					
	BLEU	METEOR	ChrF	COMET	COMET-QE	BLEURT	BLEU	METEOR	ChrF	COMET	COMET-QE	BLEURT
High-Quality Subset Baseline (HQSB)	35.45	62.00	62.75	85.50	42.00	75.90	35.62	59.90	61.11	80.50	13.50	68.10
<i>BLEU</i>												
HQSB + RL	<u>35.55</u>	62.10	62.77	86.80	<u>45.00</u>	<u>76.10</u>	36.26	60.40	61.51	<u>82.10</u>	<u>17.50</u>	<u>67.70</u>
HQSB + MBR	35.53	62.30	62.80	86.70	44.20	75.90	35.73	60.40	61.42	81.60	15.60	67.20
HQSB + RL + MBR	35.22	61.90	62.62	86.20	43.10	75.50	36.72	60.80	61.89	82.00	16.30	67.20
<i>COMET</i>												
HQSB + RL	<u>35.90</u>	<u>62.20</u>	<u>63.06</u>	86.70	44.10	75.70	<u>36.62</u>	<u>60.60</u>	<u>61.79</u>	82.20	17.40	67.60
HQSB + MBR	33.58	60.70	61.48	88.00	<u>47.90</u>	76.50	34.89	59.60	60.94	85.00	<u>27.00</u>	<u>69.80</u>
HQSB+ RL + MBR	34.92	61.80	62.84	<u>88.10</u>	47.60	<u>76.90</u>	35.97	60.20	61.45	84.40	24.50	69.20
<i>COMET-QE</i>												
HQSB + RL	35.96	62.30	63.07	86.70	44.70	75.90	<u>36.25</u>	<u>60.50</u>	<u>61.58</u>	82.40	17.70	68.10
HQSB + N-RR	31.46	58.70	60.41	87.10	53.80	75.90	29.99	54.80	56.87	82.80	39.10	66.20
HQSB + RL + N-RR	32.73	59.80	61.32	87.30	53.20	76.30	32.61	57.40	58.96	83.40	36.10	67.60
HQSB + N-RR + MBR w/ COMET	33.73	60.90	61.79	88.10	49.60	76.70	34.34	59.40	60.69	84.80	29.40	69.50
HQSB + RL + MBR w/ COMET	34.61	61.60	62.72	88.20	50.10	77.20	35.47	59.90	61.26	<u>84.90</u>	28.80	70.00
MODEL	IWSLT2017 EN→DE						IWSLT2017 EN→FR					
	BLEU	METEOR	ChrF	COMET	COMET-QE	BLEURT	BLEU	METEOR	ChrF	COMET	COMET-QE	BLEURT
Normal Baseline (NB)	32.75	62.40	60.04	84.80	38.30	74.80	41.47	68.40	66.20	84.40	21.70	73.30
<i>BLEU</i>												
NB + RL	34.48	62.90	60.51	<u>85.20</u>	<u>39.70</u>	74.40	44.58	68.60	<u>66.76</u>	<u>85.20</u>	<u>24.70</u>	72.70
NB + MBR	33.87	62.20	60.05	85.00	38.90	<u>74.50</u>	44.08	68.70	66.52	<u>85.20</u>	24.40	<u>73.20</u>
NB + RL + MBR	34.46	62.50	60.22	85.00	39.00	74.10	44.25	68.30	66.50	85.00	24.20	72.40
<i>COMET</i>												
NB + RL	<u>34.17</u>	<u>62.20</u>	59.88	85.10	39.30	74.40	<u>44.48</u>	68.70	<u>66.74</u>	85.20	24.60	72.80
NB + MBR	33.33	62.10	<u>59.97</u>	86.70	43.80	<u>75.60</u>	39.04	65.30	63.32	<u>86.80</u>	<u>37.40</u>	<u>75.00</u>
NB + RL + MBR MBR	33.75	61.90	59.72	86.10	41.80	74.90	44.24	68.50	66.62	86.30	28.30	73.60
<i>COMET-QE</i>												
NB + RL	34.53	62.90	<u>60.49</u>	85.30	40.00	74.70	<u>44.56</u>	68.70	66.87	85.30	24.90	72.90
NB + N-RR	32.31	60.70	59.06	86.40	50.00	75.60	42.48	67.20	65.38	86.60	38.30	74.00
NB + RL + N-RR	32.98	61.50	59.48	86.40	48.70	75.40	43.29	67.50	65.90	86.50	36.00	73.70
NB + N-RR + MBR w/ COMET	33.53	61.90	59.95	86.70	46.00	75.80	39.41	65.40	63.42	87.00	40.00	75.30
NB + RL + MBR w/ COMET	34.18	62.50	60.27	86.60	43.50	75.40	44.07	68.20	66.55	86.70	32.50	74.00

Table 2: Automatic evaluation metrics for the best baseline in each dataset and its variations with RL training, reranking (N -RR) and MBR decoding. BLEU, COMET, and COMET-QE serve as reward models in the context of RL training and are subjected to comparison with respect to both reranking strategies employed as the optimization metric (reranker). Best-performing values are **bolded** and best for each specific group are underlined.

intensive than MLE training), involving iterative fine-tuning of NMT models, making it suitable for capturing nuanced quality improvements from the reward models. In contrast, MBR decoding, focused on optimizing translation during inference, requires recomputation for each input sentence, allowing for computational efficiency when performed infrequently. However, it may not fully utilize the capabilities of the NMT model and can be computationally demanding in high-throughput scenarios. The choice between RL training and MBR decoding depends on specific MT system requirements, considering computational resources, translation quality objectives, and the need for real-time adaptability.

In summary, the results demonstrate that integrating RL training consistently improves translation quality in both EN→DE and EN→FR tasks across various metrics. It consistently outperforms the MLE baseline and is superior in lexical metrics scores compared to reranking strategies, which perform well according to COMET and COMET-QE.

Additionally, most top-performing models incorporate RL training, highlighting its effectiveness in complementing reranking strategies to further improve translation quality.

5 Related Work

RL-based NMT. Extensive research has been conducted on RL algorithms to improve MT. Studies by Wu et al. (2018) and Kiegedland and Kreutzer (2021) have explored the impact of RL training on large-scale translation tasks and demonstrated the effectiveness of policy gradient algorithms in mitigating exposure bias and optimizing beam search in NMT. However, both studies were limited to the use of BLEU as a reward model. Our research differs in that we explore the benefits of employing more robust preference models to improve translation quality. Additionally, other researchers have made progress in advancing reward-aware training methods. For instance, Donato et al. (2022) introduced a distributed policy gradient algorithm using mean

	WMT16 EN→DE		WMT15 EN→FR		IWSLT2017 EN→DE		IWSLT2017 EN→FR	
Method	Training	Inference	Training	Inference	Training	Inference	Training	Inference
MLE	480	5	373	3	1020	13	905	16
RL	288	5	242	3	354	13	403	16
MBR	0	212	0	55	0	500	0	660
<i>N</i> -RR	0	183	0	50	0	455	0	625

Table 3: Wall-clock time values, in minutes, that represent the efficiency of MLE, RL, MBR decoding and *N*-best reranking. The training was performed on the WMT16 EN→DE and WMT15 EN→FR high-quality subsets and on IWSLT2017 EN→DE and EN→FR entire datasets with 500 000, 300 000, 215 000 and 242 000 sentence pairs, respectively. The inference was conducted on WMT16 EN→DE, WMT15 EN→FR, IWSLT2017 EN→DE and IWSLT2017 EN→FR official test set partitions with 2999, 1500, 8079 and 8597 sentence pairs, respectively. This assessment was done with COMET as the reward model for RL and as a reranker for the reranking methods.

absolute deviation (MAD) for improved training, excelling with BLEU rewards and generalizing well to other metrics. Moreover, Ouyang et al. (2022) pioneered reinforcement learning from human feedback (RLHF) for a human-based reward model, while Gulcehre et al. (2023) proposed Reinforced Self-Training (ReST) for more efficient translation quality improvement using offline RL algorithms.

Reranking methods for NMT. Shen et al. (2004) initially introduced the concept of discriminative reranking for Statistical Machine Translation, which was later adopted by Lee et al. (2021) to train a NMT model through a reranking strategy based on BLEU. Extending this concept, MBR decoding (Kumar and Byrne, 2004) has regained popularity for candidate generation during decoding, with Müller and Sennrich (2021) finding it more robust than MAP decoding, mitigating issues like hallucinations. Furthermore, Freitag et al. (2022a) showed that coupling MBR with BLEURT, a neural metric, enhances human evaluation results when compared to lexical metrics. Fernandes et al. (2022) conducted a comprehensive study comparing various reranking strategies, including reranking and MBR decoding, with both reference-based and quality estimation metrics, concluding that these strategies lead to better translations despite the increased computational cost. In our work, we build on these foundations and show that reranking methods can be coupled with RL training to provide better translation quality to MT systems.

Data filtering for NMT. In their study, Taghipour et al. (2011) explored the use of outlier detection techniques to refine parallel corpora for MT. Meanwhile, Cui et al. (2013) proposed an unsupervised method to clean bilingual data using a random walk

algorithm that computes the importance quality score of each sentence pair and selects the higher scores. Xu and Koehn (2017) presented the Zipporah system, which is designed to efficiently clean noisy web-crawled parallel corpora. Carpuat et al. (2017) focused on identifying semantic differences between sentence pairs using a cross-lingual textual entailment system. Wang et al. (2018) proposed an online denoising approach for NMT training by using trusted data to help models measure noise in sentence pairs. Artetxe and Schwenk (2019) introduced LASER based on a BiLSTM encoder that can handle 93 different languages. Our work builds on these previous studies as we implement a data filtering method based on COMET-QE, a preference model trained on human preferences. Our approach is similar to that of Bane and Zaretskaya (2021) but is significantly more robust as preference models are much more closely aligned to human judgments compared to cross-lingual encoders.

6 Conclusion

Our thorough analysis of feedback integration methods underscores the importance of meticulous data curation for enhancing MT reliability and efficiency. Our findings demonstrate the consistent improvement in translation quality when employing neural metrics, such as COMET(-QE), during training and/or inference. RL training with data filtering stands out as significantly superior to both MLE and reranking methods. Additionally, coupling RL training with reranking techniques can further enhance translation quality. While computational efficiency remains a concern due to the added overhead of RL and reranking methods on top of MLE-trained models, their adoption should be tailored to specific task and environmental requirements.

Acknowledgments

This work was supported by EU’s Horizon Europe Research and Innovation Actions (UTTER, contract 101070631), by the project DECOLLAGE (ERC-2022-CoG 101088763), by the Portuguese Recovery and Resilience Plan through project C645008882-00000055 (Center for Responsible AI), and by Fundação para a Ciência e Tecnologia through contract UIDB/50008/2020.

References

- Artetxe, Mikel and Holger Schwenk. 2019. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *Transactions of the Association for Computational Linguistics*, 7:597–610, November.
- Bahdanau, Dzmitry, Kyung Hyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*.
- Bane, Fred and Anna Zaretskaya. 2021. Selecting the best data filtering method for NMT training. In *Proceedings of Machine Translation Summit XVIII: Users and Providers Track*, pages 89–97, Virtual, August. Association for Machine Translation in the Americas.
- Banerjee, Satanjeev and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan, June. Association for Computational Linguistics.
- Bengio, Samy, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. 2015. Scheduled sampling for sequence prediction with recurrent neural networks. *Advances in neural information processing systems*, 28.
- Bhattacharyya, Sumanta, Amirmohammad Rooshenas, Subhajit Naskar, Simeng Sun, Mohit Iyyer, and Andrew McCallum. 2021. Energy-based reranking: Improving neural machine translation using energy-based models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4528–4537, Online, August. Association for Computational Linguistics.
- Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Barry Haddow, Matthias Huck, Chris Hokamp, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Carolina Scarton, Lucia Specia, and Marco Turchi. 2015. Findings of the 2015 workshop on statistical machine translation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, Lisbon, Portugal, September. Association for Computational Linguistics.
- Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. Findings of the 2016 conference on machine translation. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 131–198, Berlin, Germany, August. Association for Computational Linguistics.
- Carpuat, Marine, Yogarshi Vyas, and Xing Niu. 2017. Detecting cross-lingual semantic divergence for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 69–79, Vancouver, August. Association for Computational Linguistics.
- Castricato, Louis, Alex Havrilla, Shahbuland Matiana, Duy V. Phung, Aman Tiwari, Jonathan Tow, and Maksym Zhuravinsky. 2023. trlx: A scalable framework for rlhf, jun.
- Cettolo, Mauro, Christian Girardi, and Marcello Federico. 2012. WIT3: Web inventory of transcribed and translated talks. In *Proceedings of the 16th Annual Conference of the European Association for Machine Translation*, pages 261–268, Trento, Italy, May 28–30. European Association for Machine Translation.
- Cettolo, Mauro, Marcello Federico, Luisa Bentivogli, Jan Niehues, Sebastian Stüker, Katsuhito Sudoh, Koichiro Yoshino, and Christian Federmann. 2017. Overview of the IWSLT 2017 evaluation campaign. In *Proceedings of the 14th International Conference on Spoken Language Translation*, pages 2–14, Tokyo, Japan, December 14–15. International Workshop on Spoken Language Translation.
- Conneau, Alexis, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.
- Cui, Lei, Dongdong Zhang, Shujie Liu, Mu Li, and Ming Zhou. 2013. Bilingual data cleaning for SMT using graph-based random walk. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 340–345, Sofia, Bulgaria, August. Association for Computational Linguistics.
- Deutsch, Daniel, Rotem Dror, and Dan Roth. 2022. On the limitations of reference-free evaluations of generated text. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10960–10977, Abu Dhabi, United Arab

- Emirates, December. Association for Computational Linguistics.
- Donato, Domenic, Lei Yu, Wang Ling, and Chris Dyer. 2022. Mad for robust reinforcement learning in machine translation.
- Eikema, Bryan and Wilker Aziz. 2022. Sampling-based approximations to minimum bayes risk decoding for neural machine translation.
- Fan, Angela, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia, July. Association for Computational Linguistics.
- Fernandes, Patrick, António Farinhas, Ricardo Rei, José GC de Souza, Perez Ogayo, Graham Neubig, and André FT Martins. 2022. Quality-aware decoding for neural machine translation. *arXiv preprint arXiv:2205.00978*.
- Freitag, Markus, David Grangier, Qijun Tan, and Bowen Liang. 2022a. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.
- Freitag, Markus, Ricardo Rei, Nitika Mathur, Chikui Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022b. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Gulcehre, Caglar, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, Wolfgang Macherey, Arnaud Doucet, Orhan Firat, and Nando de Freitas. 2023. Reinforced self-training (rest) for language modeling.
- Holtzman, Ari, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration.
- Jiao, Wenxiang, Xing Wang, Shilin He, Irwin King, Michael R. Lyu, and Zhaopeng Tu. 2020. Data rejuvenation: Exploiting inactive training examples for neural machine translation.
- Khayrallah, Huda and Philipp Koehn. 2018. On the impact of various types of noise on neural machine translation. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 74–83, Melbourne, Australia, July. Association for Computational Linguistics.
- Kiegeland, Samuel and Julia Kreutzer. 2021. Revisiting the weaknesses of reinforcement learning for neural machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1673–1681.
- Kingma, Diederik P. and Jimmy Ba. 2017. Adam: A method for stochastic optimization.
- Koehn, Philipp and Rebecca Knowles. 2017. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver, August. Association for Computational Linguistics.
- Koehn, Philipp, Vishrav Chaudhary, Ahmed El-Kishky, Naman Goyal, Peng-Jen Chen, and Francisco Guzmán. 2020. Findings of the WMT 2020 shared task on parallel corpus filtering and alignment. In Barraud, Loïc, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Yvette Graham, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, and Matteo Negri, editors, *Proceedings of the Fifth Conference on Machine Translation*, pages 726–742, Online, November. Association for Computational Linguistics.
- Kong, Xiang, Zhaopeng Tu, Shuming Shi, Eduard Hovy, and Tong Zhang. 2018. Neural machine translation with adequacy-oriented learning.
- Kreutzer, Julia, Joshua Uyheng, and Stefan Riezler. 2018. Reliability and learnability of human bandit feedback for sequence-to-sequence reinforcement learning.
- Kumar, Shankar and William Byrne. 2004. Minimum Bayes-risk decoding for statistical machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 169–176, Boston, Massachusetts, USA, May 2 - May 7. Association for Computational Linguistics.
- Lample, Guillaume, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. 2017. Unsupervised machine translation using monolingual corpora only. *arXiv preprint arXiv:1711.00043*.
- Lee, Ann, Michael Auli, and Marc’Aurelio Ranzato. 2021. Discriminative reranking for neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7250–7264.
- Malli, Marilena and George Tambouratzis. 2022. Evaluating corpus cleanup methods in the WMT’22 news translation task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 335–341, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.

- Müller, Mathias and Rico Sennrich. 2021. Understanding the properties of minimum Bayes risk decoding in neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 259–272, Online, August. Association for Computational Linguistics.
- Ng, Nathan, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, and Sergey Edunov. 2019. Facebook FAIR’s WMT19 news translation task submission. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 314–319, Florence, Italy, August. Association for Computational Linguistics.
- Och, Franz Josef. 2003. Minimum error rate training in statistical machine translation. In *Proceedings of the 41st annual meeting of the Association for Computational Linguistics*, pages 160–167.
- Ott, Myle, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation. In Dy, Jennifer and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3956–3965. PMLR, 10–15 Jul.
- Ouyang, Long, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Perez, Ethan, Saffron Huang, H. Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. Red teaming language models with language models. *CoRR*, abs/2202.03286.
- Popović, Maja. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal, September. Association for Computational Linguistics.
- Puterman, Martin L. 1990. Markov decision processes. *Handbooks in operations research and management science*, 2:331–434.
- Raffel, Colin, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *CoRR*, abs/1910.10683.
- Ranzato, Marc’Aurelio, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2016. Sequence level training with recurrent neural networks. In *4th International Conference on Learning Representations, ICLR 2016*.
- Reddy, Raj. 1977. Speech understanding systems: A summary of results of the five-year research effort at carnegie mellon university.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Unbabel’s participation in the WMT20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 911–920, Online, November. Association for Computational Linguistics.
- Rei, Ricardo, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Rei, Ricardo, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Schulman, John, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347.
- Sellam, Thibault, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online, July. Association for Computational Linguistics.
- Shen, Libin, Anoop Sarkar, and Franz Josef Och. 2004. Discriminative reranking for machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 177–184.
- Shen, Shiqi, Yong Cheng, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. 2016. Minimum risk training for neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1683–1692.
- Stiennon, Nisan, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,

- Dario Amodei, and Paul Christiano. 2022. Learning to summarize from human feedback.
- Sutton, Richard S. and Andrew G. Barto. 2018. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA.
- Taghipour, Kaveh, Shahram Khadivi, and Jia Xu. 2011. Parallel corpus refinement as an outlier detection algorithm. In *Proceedings of Machine Translation Summit XIII: Papers*, Xiamen, China, September 19-23.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Wang, Wei, Taro Watanabe, Macduff Hughes, Tetsuji Nakagawa, and Ciprian Chelba. 2018. Denoising neural machine translation training with trusted data and online data selection.
- Williams, Ronald J. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256.
- Wiseman, Sam and Alexander M Rush. 2016. Sequence-to-sequence learning as beam-search optimization. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1296–1306.
- Wolf, Thomas, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Huggingface’s transformers: State-of-the-art natural language processing.
- Wu, Lijun, Fei Tian, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. 2018. A study of reinforcement learning for neural machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3612–3621.
- Xu, Hainan and Philipp Koehn. 2017. Zipporah: a fast and scalable data cleaning system for noisy web-crawled parallel corpora. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2945–2950, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Zoph, Barret, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation.

A Additional Results

To gain deeper insights into the effectiveness of both training and inference techniques, we also conducted a small fine-grained study evaluating the translation quality of models. Specifically, we compared translations produced by the High-Quality Subset Baseline using three different methods: MBR with COMET, RL training with COMET-QE as a reward model and a hybrid approach combining both. This complementary evaluation primarily relies on BLEURT, a neural metric highly correlated with human judgments and independent from the used reward models.

The overall BLEURT scores for these systems can be obtained from Table 2, with HQSB, HQSB + MBR w/ COMET, HQSB + RL w/ COMET-QE and HQSB + RL w/ COMET-QE + MBR w/ COMET having 75.90, 76.50, 75.94 and 77.20, respectively. Figure 4 illustrates a discernible trend: across varying lengths of source sentences, the model trained with RL and employing MBR during inference consistently yields translations of higher quality. Additionally, there is a noticeable decline in translation quality when MBR alone is employed for exceptionally long sentences, a phenomenon seemingly linked to specific hallucinations evident in Figure 3. Furthermore, Table 4 showcases the most critical examples of hallucinations obtained during this analysis.

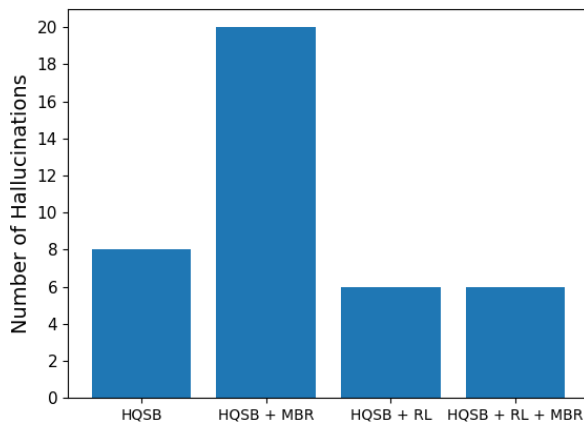


Figure 3: Number of hallucinations on the WMT16 EN→DE test set with 3000 sentences.

Examining Figures 5 and 6, depicting sentence counts across various ranges of BLEU and BLEURT scores, respectively, reveals the trend that the HQSB + RL + MBR system consistently outperforms the remaining systems across both metrics. Once again, the prevalence of low BLEU scores

Source: Posted by TODAY on Monday, September 14, 2015

Reference: Geschrieben von TODAY am Montag, 14. September 2015

MBR Hallucination: Posted by TODAY am Montag, 14. September 2015, 14:45 Uhr Posted by TODAY am Montag, September 14, 2015, 14:40 Uhr Posted by TODAY am Montag, September 14, 2015, 14:00 Uhr Posted by TODAY am Montag, September 14, 2015, 14:30 Uhr Posted by TODAY am Montag, September 14, 2015, 14:30 Uhr Posted by TODAY am Montag, September 14, 2015, 14:30 Uhr Posted by TO

RL + MBR Translation: Veröffentlicht von TODAY am Montag, 14. September 2015

Source: Seehofer: "Borders will not be cordoned off"

Reference: Seehofer: "Grenzen werden nicht abgeriegelt"

MBR Hallucination: Seehofer: "Grenzen werden nicht abgeschottet" Seehofer: "Grenzen werden nicht abgeschottet" Seehofer: "Grenzen werden nicht abgeschottet" Seehofer: "Grenzen werden nicht abgeschottet" Seehofer: "Grenzen werden nicht abgeschottet" Seehofer: "Grenzen werden nicht abgeschottet" Seehofer: "Grenzen werden nicht abgeschottet" Seehofer: "Grenzen werden nicht abgeschottet" Seehofer

RL + MBR Translation: Seehofer: "Grenzen werden nicht abgeriegelt"

Source: Croatia: "We are letting the refugees through"

Reference: Kroatien: "Wir lassen die Flüchtlinge durch"

MBR Hallucination: Kroatien: "Wir lassen die Flüchtlinge durch" "Wir lassen die Flüchtlinge durch" Kroatien: "Wir lassen die Flüchtlinge durch" Kroatien: "Wir lassen die Flüchtlinge durch" Kroatien: "Wir lassen die Flüchtlinge durch" Kroatien: "Wir lassen die Flüchtlinge durch" Kroatien: "Wir lassen die Flüchtlinge durch"

RL + MBR Translation: Kroatien: "Wir lassen die Flüchtlinge durch"

Table 4: Instances of oscillatory hallucinations generated by the HQSB + MBR model.

underscores the issue of hallucinations associated with MBR. Furthermore, HQSB and HQSB + RL systems are quite competitive but a slight edge must be given to RL in enhancing the performance of the models

The bucketed word accuracy analysis aims to evaluate how effectively each system is at generating different types of words. Figure 7 shows that all four systems demonstrate robustness across all word frequencies but perform significantly better with higher-frequency words. Notably, among these systems, the one integrating reinforcement learning (RL) emerges as the top performer, emphasizing its effectiveness in word generation tasks.

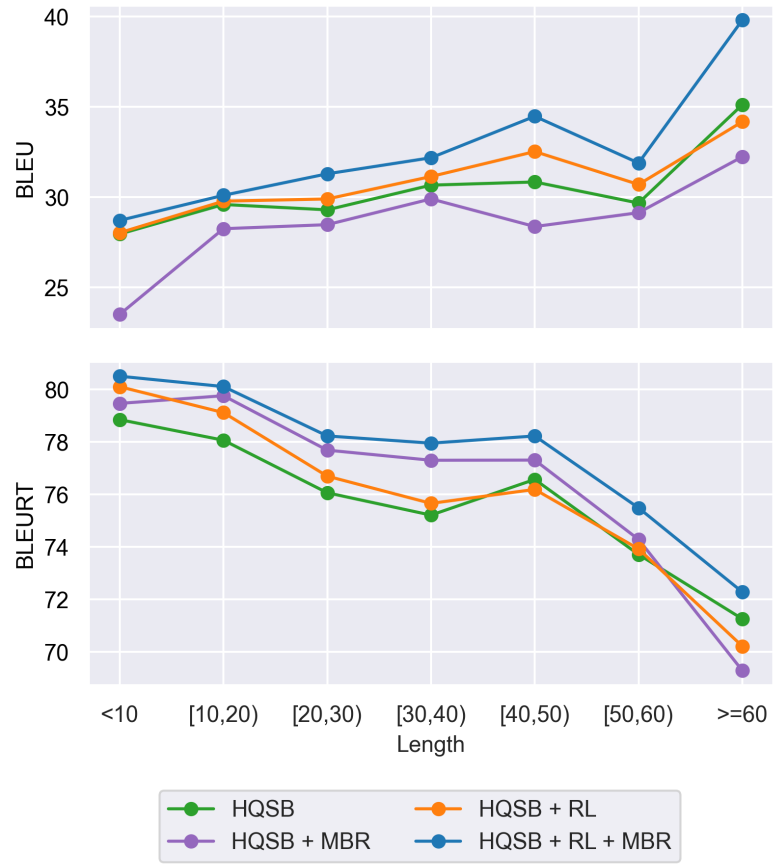


Figure 4: Comparison of BLEU (top) and BLEURT (bottom) scores for WMT16 EN→DE translations across diverse source sentence lengths, highlighting the influence of sentence length on translation quality.

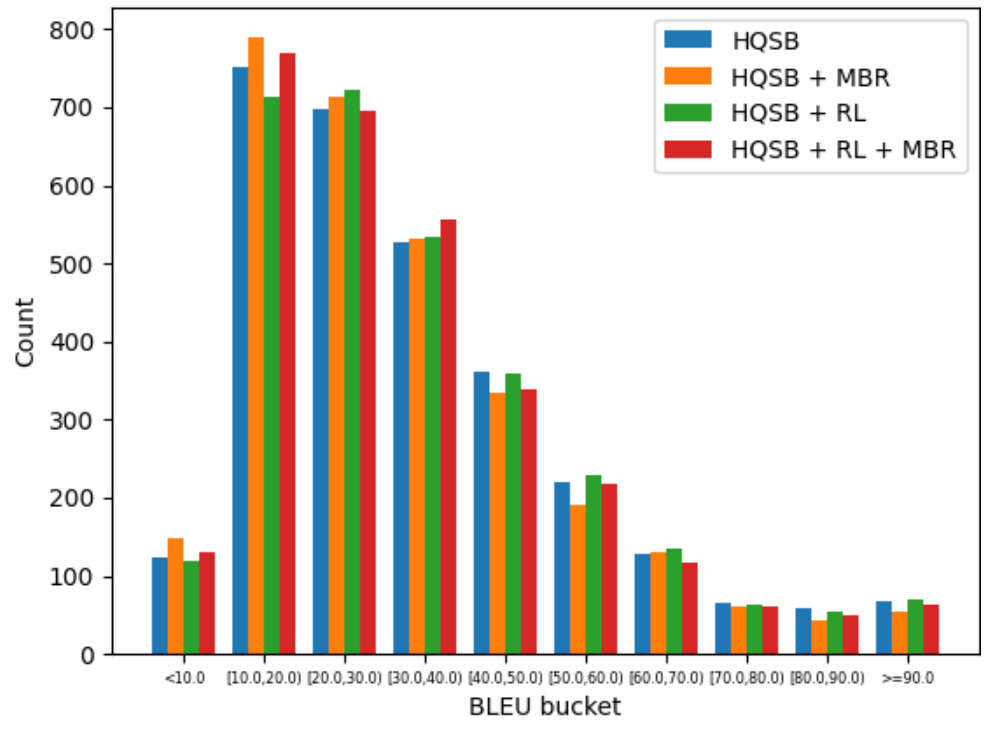


Figure 5: Histograms of sentence BLEU scores for the specific systems on WMT16 EN→DE.

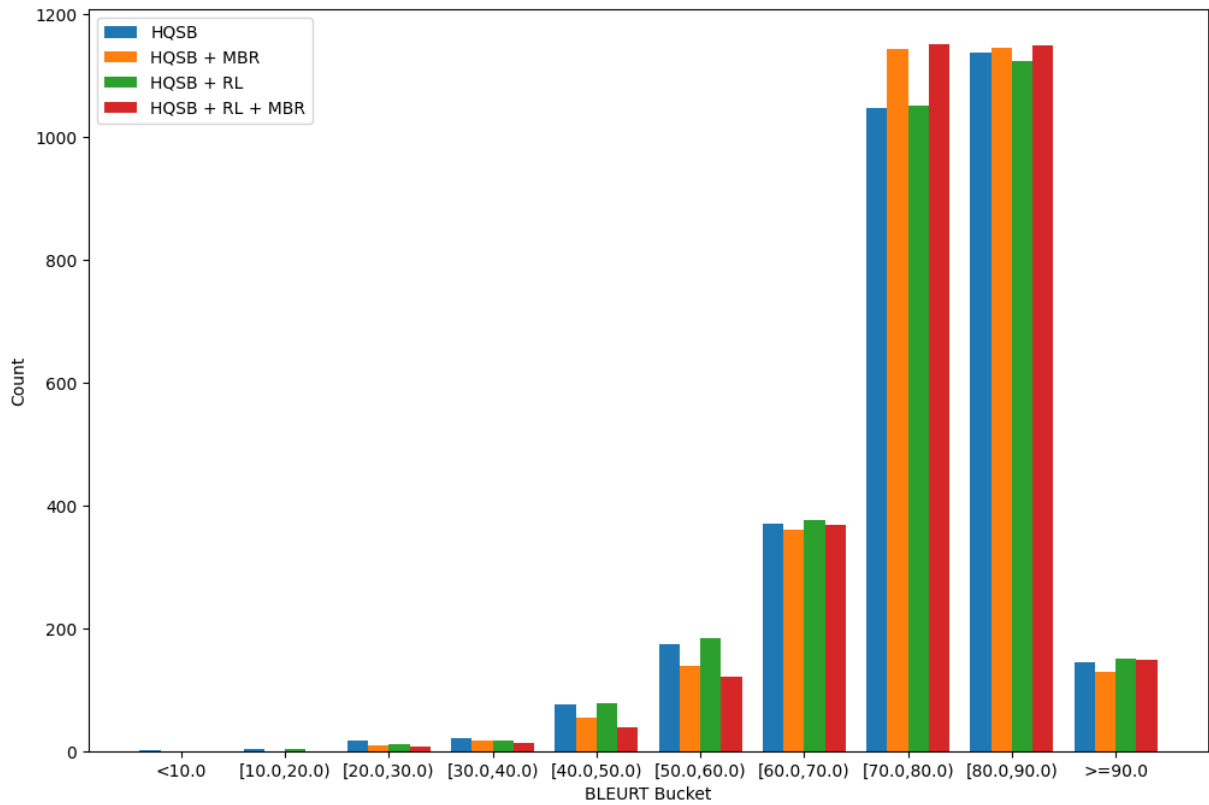


Figure 6: Histograms of sentence BLEURT scores for the specific systems on WMT16 EN→DE.

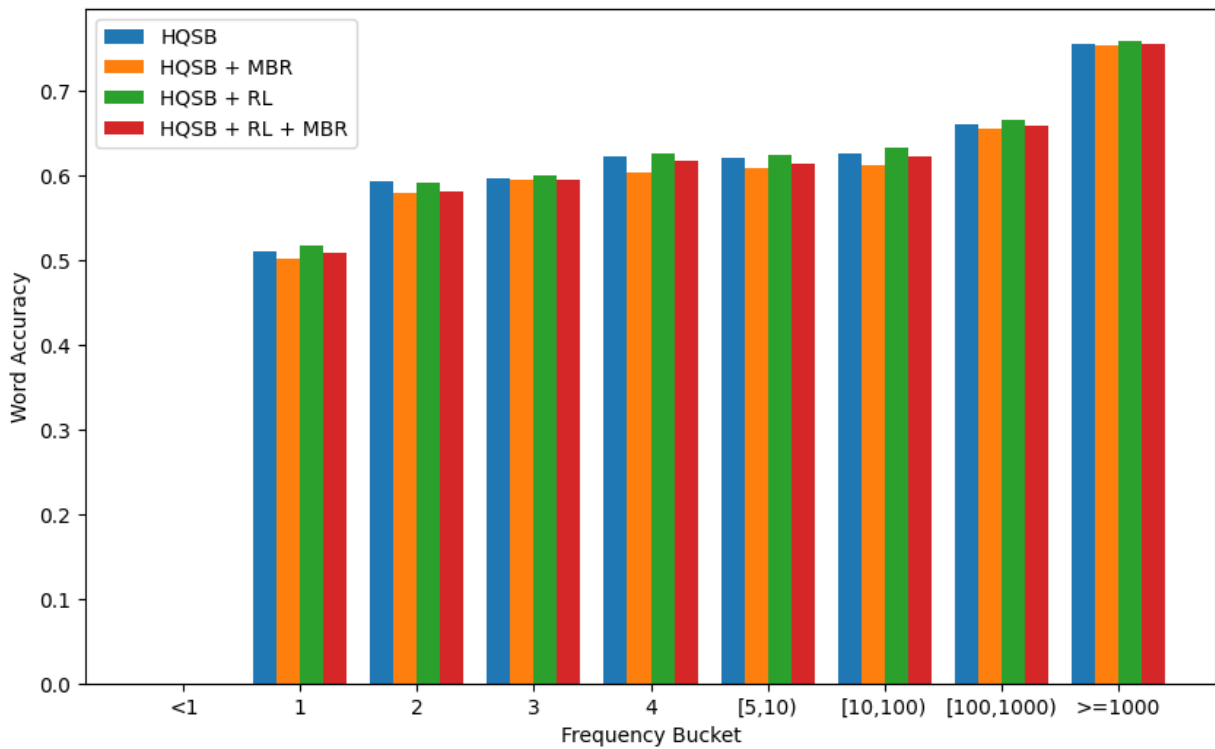


Figure 7: Word F-Measure Bucketed by Frequency for the specific systems on WMT16 EN→DE.

Enhancing Scientific Discourse: Machine Translation for the Scientific Domain

Dimitris Roussis, Sokratis Sofianopoulos, Stelios Piperidis

Institute for Speech and Language Processing, Athena RC
Artemidos 6 & Epidavrou, Athens, Greece
{dimitris.roussis, s_sofian, spip}@athenarc.gr

Abstract

The increasing volume of scientific research necessitates effective communication across language barriers. Machine translation (MT) offers a promising solution for accessing international publications. However, the scientific domain presents unique challenges due to its specialized vocabulary and complex sentence structures. In this paper, we present the development of a collection of parallel and monolingual corpora for the scientific domain. The corpora target the language pairs Spanish-English, French-English, and Portuguese-English. For each language pair, we create a large general scientific corpus as well as four smaller corpora focused on the domains of: Cancer Research, Energy Research, Neuroscience, and Transportation research. To evaluate the quality of these corpora, we utilize them for fine-tuning general-purpose neural machine translation (NMT) systems. We provide details regarding the corpus creation process, the fine-tuning strategies employed, and we conclude with the evaluation results.

1 Introduction

The growth of scientific research across disciplines has intensified the need for efficient communication and international collaboration that transcends language barriers. While English is the dominant language of scientific publications (Altbach,

2007), a substantial volume of valuable academic work is produced in other languages. According to a 2019 study (Stockemer, 2019), almost 40% of the articles of non-anglophone researchers are submitted in a language other than English. (Altbach, 2007). Machine translation (MT) offers a compelling solution, providing access to a vast pool of international research and fostering seamless collaboration among researchers worldwide.

Advancements in Neural Machine Translation have significantly improved the quality of translations in various domains, largely driven by the application of the Transformer architecture (Vaswani et al., 2017). Transformers revolutionized NMT by enabling efficient parallel processing of entire sequences, leading to significant improvements in translation quality and fluency. However, the performance of NMT models often suffers when translating specialized domains due to the presence of specific terminology and sentence structures. Translating scientific text presents unique challenges distinct from general language translation (Byrne, 2012). Scientific domains are characterized by:

- **Specialized Lexicon:** These domains employ a rich vocabulary of technical terms and abbreviations often absent from general language corpora.
- **Syntactic Complexity:** Scientific writing frequently utilizes complex sentence structures to convey precise relationships and subtle meanings.
- **Domain-Specific Discourse:** Each scientific domain possesses its own unique discourse patterns and conventions that are critical to understand for accurate translation.

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

- **Indeterminacy:** The division of scientific areas into numerous sub-disciplines and the existence of multidisciplinary works further complicate the definition of what a scientific domain encompasses.

These factors can significantly impair the performance of generic NMT models, leading to mistranslations or loss of scientific meaning.

To address these challenges, this paper presents the development of domain-specific corpora for scientific research and their application towards creating NMT models for four academic domains. We aimed at using the open-source NMT models which exhibit the best generic performance as base models and fine-tune them on these domain-specific corpora, so as to achieve translations that are not only fluent but also accurate and faithful to the scientific content. We focus on the Spanish–English (ES–EN), French–English (FR–EN), and Portuguese–English (PT–EN) translation directions, creating parallel corpora for each, alongside monolingual corpora for the respective languages (i.e., English, Spanish, French, and Portuguese). Our corpora encompass a large general scientific corpus and smaller sub-corpora dedicated to the research areas of:

- Cancer Research
- Energy Research
- Neuroscience
- Transportation Research

We begin by presenting our dataset creation process which exploits the wealth of parallel titles and abstracts from bachelor and master theses, doctoral dissertations, and other scientific publications (such as published books and articles). Next, we outline the methodology of fine-tuning pre-trained NMT models using the aforementioned datasets. Finally, we evaluate the resulting models’ performance by contrasting their output with the translations received from the original general-purpose NMT models as well as Google Translate.

2 Related Work

Even though there is a plethora of parallel corpora in the ever-growing OPUS collection (Tiedemann, 2012), there is a shortage of those which aim at addressing parallel data acquisition for the scientific and academic domains, relatively to the importance and challenges of scientific translation.

In OPUS, there are 2 parallel datasets of note, namely CAPES (Soares et al., 2018) and SciELO (Soares et al., 2019). CAPES originates from the database of the Theses and Dissertations Catalog (TDC) and contains approximately 1.2 million sentence pairs for the EN–PT language pair (Soares et al., 2018), mined from theses and dissertation abstracts of students in post-graduate programs across Brazilian universities. Similarly, the SciELO parallel corpus has been extracted through the use of the SciELO database, which contains a broad range of open-access scientific articles. SciELO consists of approximately 3.3 million parallel sentences and metadata for English, Spanish, and Portuguese, some of which are trilingually aligned (i.e., EN–ES–PT) (Soares et al., 2019). Both corpora are evaluated manually, as well as automatically by training and evaluating MT systems.

In order to address the parallel data gap for scientific texts among underrepresented European languages, the SciPar corpus (Roussis et al., 2022) was created and made publicly available via the ELRC-SHARE repository.¹ SciPar contains 9.17 million sentence pairs in 31 language pairs and has been constructed from the titles and abstracts of bachelor and master theses, doctoral dissertations, and other scientific publications. It has been constructed through collecting, parsing, and processing metadata from 86 institutional repositories, digital libraries of universities, and national archives.

The *Translations and Open Science* project focused on building EN–FR NMT systems for three pilot domains: (a) Climatology and Climate Change, (b) Neurosciences, and (c) Human Mobility, Environment, and Space. Fiorini et al. (2023) collected 316,701 parallel segments and 1,112 bilingual terms for these three domains, trained a generic NMT model from scratch, and the fine-tuned it on the collected datasets. The paper describes the manual and automatic evaluation and comparison that was conducted to determine domain-specific translation quality, while it also shows that adding data from SciPar results in further improvements.

3 Dataset Creation

In this section, we outline the end-to-end pipeline used for mining high-quality parallel and monolingual corpora from the titles and abstracts of the-

¹<https://elrc-share.eu/>

Domain	EN-ES	EN-PT	EN-FR
Cancer	57,226	123,357	49,112
Energy	107,710	205,662	87,918
Neuroscience	40,467	85,717	45,650
Transp/tion	26,795	35,181	19,151
Gen. Scient.	3,913,214	5,255,552	1,648,200
Total	4,145,412	5,705,469	1,850,031

Table 1: Parallel Corpora Sizes per Language Pair and Domain

ses and dissertations. First, we detail the strategy used to process approximately 9.3 million records from 62 academic repositories to extract 11,700,912 sentence pairs. A detailed list of the repositories can be seen at Table 4 of Appendix A. Then, we present the sizes of the parallel and monolingual corpora which resulted after applying various filtering methods, and, finally, we provide a brief documentation of the process used to create benchmark developer and test sets.

3.1 Repository Processing

Our strategy is directed toward structured metadata extraction and processing of academic records in order to mine domain-specific monolingual and parallel sentences, while also facilitating differences among various repositories. It constitutes a unified framework that builds upon and extends the approach used in SciPar (Roussis et al., 2022). The various steps used in our repository processing pipeline are the following:

- **Acquisition of academic records:** Initially, we utilized the GNU Wget² package to automatically download all the records across the 62 repositories as HTML files.
- **Custom repository configuration:** In order to extract structured metadata across repositories which differ in the way they store information, we created a configuration file for each one (see Appendix B for an example). This required manual inspection of several HTML files from each repository and configure custom regex patterns and rules. Furthermore, we defined minimum character lengths for the extraction of abstracts and titles to ensure the validity of the extracted data.

²<https://www.gnu.org/software/wget/>

- **Parsing and metadata extraction:** We utilized the Beautiful Soup³ package and regex pattern processing to parse the records and extract structured metadata. This step used the custom configuration files for each repository and resulted in the construction of a JSON file for each record, organizing the data for further processing (see Appendix C for an example).
- **Domain classification:** Leveraging a simple keyword-based classification method, each record was categorized into one of four distinct domains: Cancer Research, Transportation Research, Energy Research, and Neuroscience. Records not fitting exclusively into these categories were classified as belonging to the "General Academic" domain.
- **Text extraction:** From each JSON's titles and abstracts, both monolingual and parallel documents were extracted. This process preserved information regarding the identified domain and language of each record. Candidate parallel documents were created when two or more titles and/or abstracts were present. The NLTK library⁴ was employed to split abstracts into sentences.
- **Parallel sentences mining:** Utilizing LASER and margin-based scoring (Artetxe and Schwenk, 2019a; Artetxe and Schwenk, 2019b), we extracted parallel sentences from the candidate parallel documents. We used 2 NVIDIA 2080 Ti GPUs for this process.

3.2 Domain-specific Parallel Corpora

In order to construct domain-specific parallel corpora for all the language pairs that we targeted, we

³<https://www.crummy.com/software/BeautifulSoup/>

⁴<https://www.nltk.org/>

	EN	ES	FR	PT
Cancer	393,488	76,296	6,933	33,947
Energy	342,144	228,818	9,479	66,654
Neuroscience	262,618	47,112	5,380	19,640
Transp/tion	33,509	65,374	2,934	13,376
Gen. Scient.	13,187,215	10,512,255	753,487	3,335,615
Total	14,218,974	10,929,855	778,213	3,469,232

Table 2: Monolingual Corpora Sizes per Language Pair and Domain

concatenated the parallel sentences extracted earlier from each one of the 62 repositories into domain and language pair specific files.

Additionally, a LASER alignment score threshold of 0.98 was applied which, although being lower than thresholds used in other works which mine parallel sentences in a global way, ensures that parallel data originating from titles are not discarded (Roussis et al., 2022).

Finally, we deduplicated all parallel corpora and filtered them by removing sentence pairs which: (a) have identical source and target sentences, (b) contain an empty sentence in either side, (c) consist of more than 250 words in either sentence, (d) are solely comprised of digits, (e) are identified as belonging to incorrect languages, or f consist mostly of URLs and e-mails (Papavassiliou et al., 2018).

The resulting domain-specific parallel corpora in Table 1 sum up to 11,700,912 sentence pairs in total. We can see that the the "General Scientific" domain considerably surpasses all other domains combined in corpus size across all language pairs. Furthermore, the corpus sizes for the EN-PT and EN-ES language pairs generally exceed those for the EN-FR language pair, which is indicative of the plethora of academic repositories originating from Latin American countries. Among the four focused domains, the "Energy" and "Cancer Research" domains exhibit higher data representations. The disparities in corpus sizes across different domains reflect the extent of each domain's scope, and the inherent challenges of domain classification.

In Table 2, we list the sizes of the domain-specific monolingual corpora for each of the four targeted languages. These results do not include the monolingual sentences in the parallel data and have been extracted from academic records without parallel titles and/or abstracts. It can be ob-

served that there are many variations among the domains that we focused on; The "General Scientific" domain is significantly larger than all the other domains that we focused on. However, more pronounced differences concern the data availability for each language, as, for example, the English monolingual corpus is two orders of magnitude larger than the French one. Although we did not directly use monolingual data in this work, they can potentially enhance the performance of NMT models with the use of back-translation (Sennrich et al., 2015; Edunov et al., 2018), as well as constitute a high-quality part of a pre-training dataset for a Large Language Model (LLM).

3.3 Benchmark Creation

The benchmark creation process was tailored to meet the specific demands of evaluating our fine-tuned models across the four scientific domains that we focus on. Our goal was to create:

- **Developer sets** which can be used to monitor the training process and select optimal hyper-parameters
- **Test sets** which can be used to evaluate the generalization capabilities of our models, as well as to compare them with the base OPUS-MT models and Google Translate.

For each of the four targeted domains, 1,000/1,000 sentence pairs were designated for development/testing, whereas the "General Scientific" set was allocated 3,000/3,000 sentence pairs. In order to build the domain-specific developer and test sets, we followed a structured multi-stage process. We started from the filtered parallel corpora for each domain and language pair, and applied a LASER threshold greater than 1.08 to identify sentences with substantial semantic alignment. Additionally, we used a stricter token ratio filter (1.66)

and removed sentence pairs in which either sentence is less than 3 words long.

Afterwards, we sample 2,000 unique records for each domain-specific dataset and 6,000 for the "General Scientific" datasets, from which a single sentence pair was randomly selected for inclusion in the combined developer/test sets. In other words, each parallel sentence in the sets is derived from a single thesis/abstract from a pool of randomly selected ones.

For the construction of developer and test sets, we equally distributed the sentence pairs into two, resulting in 1,000 parallel sentences for each domain-specific developer and test set and 3,000 parallel sentences for the "General Scientific" developer and test set.

4 Domain Adaptation for NMT

Using domain adaptation, we can enhance the ability of an existing MT system to translate documents from a very specific domain, through the utilization of an in-domain parallel corpus. To this end, we chose to use selected pre-trained models from OPUS-MT (Tiedemann, 2020) as our baseline systems. We chose these models because they are all based on state-of-the-art transformer-based neural machine translation architectures, and they were trained on freely available parallel corpora from the OPUS3 bitext repository. Even though these models are not considered to produce the best quality for the selected language pairs, they provide a robust initial performance across a variety of language pairs. Additionally, they have also been trained on general domain texts, making them ideal for our purpose, as they have seen little to no data from the scientific domain.

Since we decided to work on French, Spanish and Portuguese, we selected the open-source OPUS-MT models (Tiedemann and Thottingal, 2020) for the FR-EN, ES-EN, and PT-EN, translation directions as our base models. We aimed at using the largest available OPUS-MT models for each language pair and made use of the latest ones which were trained for the Tatoeba Challenge (Tiedemann, 2020), an open project encouraging people to develop machine translation in real-world cases for many languages. In particular, we used the following models:

- **FR-EN:** We selected the transformer-big variant, trained on English-French parallel

data.⁵ We use a learning rate of $7e-5$ and 0.2 dropout for all the fine-tuning experiments.

- **ES-EN:** We chose to use the transformer-big variant which was trained for translating Catalan (CA), Occitan (OC), and Spanish (ES) texts into English,⁶ as a larger variant was not available for the ES-EN translation direction. We use a learning rate of $7e-5$ and 0.2 dropout for all the fine-tuning experiments.
- **PT-EN:** For the PT-EN translation direction, we used the transformer-align (Bahdanau et al., 2014) variant,⁷ as there were not any variants with more parameters. We use a learning rate of $5e-5$ and 0.1 dropout for all the fine-tuning experiments.

All of these variants have been originally trained on subsets of the OPUS-MT corpus, along with additional back-translated data; a method which has been widely adopted by the research community to further improve translation quality (Sennrich et al., 2015; Edunov et al., 2018). For fine-tuning the OPUS-MT models we used the MarianMT framework (Junczys-Dowmunt et al., 2018), a decision driven by MarianMT's efficiency, flexibility, and open-source nature, making it particularly suited for research purposes.

Fine-tuning was conducted using the curated corpora which resulted from the filtering process described in section 3.2, comprising parallel sentences from academic texts across our target domains: Cancer Research, Energy Research, Neuroscience, and Transportation Research. These corpora were further enriched with general academic texts to provide a broader linguistic context, thereby mitigating the potential for overfitting to domain-specific jargon and syntax. In our domain adaptation experiments, we make use of all developer sets described in Section 3.3 (7000 sentence pairs in total) so that the fine-tuning process does

⁵opusTCv20210807+bt-transformer-big_2022-03-09
<https://github.com/Helsinki-NLP/Tatoeba-Challenge/tree/master/models/fra-eng>

⁶opusTCv20210807+bt-transformer-big_2022-03-13
<https://github.com/Helsinki-NLP/Tatoeba-Challenge/tree/master/models/cat+oci+spa-eng>

⁷opus+bt-2021-04-30
<https://github.com/Helsinki-NLP/Tatoeba-Challenge/tree/master/models/por-eng>

	Average of 4 domains			General Scientific		
	BLEU	chrF2++	COMET	BLEU	chrF2++	COMET
ES-EN						
Base OPUS-MT	49.7	70.5	69.5	51	71.7	68.9
+ FT w/ IND	50.7	71.3	70.4	52.1	72.2	68.9
+ FT w/ IND & GSC	51.9	71.7	70.9	54	73.1	71
Google Translate	51	71.1	72.2	52.6	72.4	73
$\Delta(\text{Ours} - \text{Baseline})$	+2.2	+1.2	+1.4	+3	+1.4	+2.1
PT-EN						
Base OPUS-MT	46	68.3	66.7	44.9	67.7	66.3
+ FT w/ IND	48.1	69.7	66.8	46.6	68.7	66.1
+ FT w/ IND & GSC	48.4	69.9	67.6	47.3	69.1	67.8
Google Translate	48	69.3	70.6	46.7	68.6	70.5
$\Delta(\text{Ours} - \text{Baseline})$	+2.4	+1.6	+0.9	+2.4	+1.4	+1.5
FR-EN						
Base OPUS-MT	37.6	63.6	57.5	38.4	63.3	57.2
+ FT w/ IND	39.1	65	58.9	39.4	64.5	58.4
+ FT w/ IND & GSC	40.2	65.7	59.9	40.7	65.3	59.5
Google Translate	41	65	62.2	40.9	64.4	61.6
$\Delta(\text{Ours} - \text{Baseline})$	+2.6	+2.1	+2.4	+2.3	+2	+2.3

Table 3: Results of Domain Adaptation Experiments

not continue if the BLEU or chrF scores do not improve for 5 consecutive checkpoints (every 500 steps). Batch size is automatically determined by MarianNMT to fit reserved GPU memory, while parameters are updated every 2 batches (Ott et al., 2018).

During the fine-tuning phase, particular attention was paid to the balance between domain-specific and general academic texts, ensuring that the models retained their generalizability while enhancing their domain-specific performance. Model hyperparameters were meticulously optimized through a series of experiments, with evaluation metrics such as BLEU (Papineni et al., 2002), chrF2++ (Popović, 2015; Popović, 2016; Popović, 2017), and COMET (Rei et al., 2020) used to monitor improvements in translation quality.

5 Results & Evaluation

In order to determine the improvement on translation quality for each specific domain of interest, we perform evaluation using all the 4 domain-specific test sets, as well as the General Scientific test set (see Section 3.3). To that end, the most

widely-used benchmarks for the evaluation of automatic translations have been chosen, these being BLEU, chrF2++, and COMET. BLEU measures how much overlap there is between the translated text and a reference translation in terms of n-gram phrases, while chrF2++ considers character-level accuracy and morphology in addition to n-grams. Unlike BLEU and chrF2++ which focus on word-level overlap, COMET doesn’t directly compare the words themselves. Instead, it uses the XLM-R Large as its backbone and estimates how well the translated text captures the meaning conveyed in the reference text. COMET, in particular, has been shown to be superior to the other two metrics that we used, as it better correlates with human judgements and provides a robust way to differentiate high-performing systems (Rei et al., 2020).

We use the SacreBLEU toolkit (Post, 2018) for evaluation of BLEU⁸ and chrF2++⁹ as it ensures that they can be computed in a reproducible way, while the wmt20-comet-da¹⁰ version was used for

⁸nrefs:1| bs:1000| seed:12345| case:mixed| eff:no| tok:13a| smooth:exp| version:2.0.0

⁹nrefs:1| bs:1000| seed:12345| case:mixed| eff:yes| nc:6| nw:2| space:no| version:2.0.0

¹⁰<https://huggingface.co/Unbabel/>

computing the COMET score.

In Table 2, we can see that our fine-tuned models outperform the base models. In all language pairs and across all metrics, the systems that use fine-tuning with in-domain data (either alone or combined with general scientific text) outperform the baseline OPUS-MT system by +2.4 BLEU and +1.6 COMET on average. Fine-tuning the baseline OPUS-MT models with in-domain parallel data (+FT w/IND) yields an average improvement of +1.5 BLEU and +0.8 COMET.

Across all language pairs, adding general scientific text for fine-tuning on top of in-domain data (+FT w/IND & GSC) leads to further improvements compared to just using in-domain data, improving by +0.9 BLEU and +0.8 COMET. This suggests that general scientific data can act as a helpful supplement for domain adaptation, although it should be noted that their number is an order of magnitude larger than the in-domain data.

While the fine-tuned systems achieve the best scores in the majority of cases, Google Translate remains highly competitive, producing the highest score in COMET and outperforming all other systems in almost all metrics for FR–EN. This indicates that even without specific domain adaptation, Google Translate offers strong translation capabilities for high-resource language pairs, even for niche domains like the ones selected for this experiment.

6 Conclusion

Our work presents the development of domain-specific corpora for scientific research and their application towards creating improved NMT models for four scientific domains: Cancer Research, Energy Research, Neuroscience, and Transportation Research. We acquired 11.7M parallel sentences for three language pairs (EN–ES, EN–PT, and EN–FR) by processing the records from 62 academic repositories. Our aim was to leverage the strengths of existing open-source NMT models while specializing them in these particular domains through fine-tuning.

Our findings indicate that fine-tuning generic NMT models with domain-specific parallel data leads to substantial improvements in translation quality for the targeted scientific domains. Additionally, by including general scientific text alongside domain-specific data during fine-tuning of-

fers further enhancements, potentially by providing broader linguistic context.

This work contributes to the ongoing effort to bridge the language gap in scientific research by developing NMT models that can accurately and fluently translate scientific text across various domains.

Looking ahead, there are a lot of ideas that we could implement in order to further improve our results. Multi-domain adaptation holds promise for creating models adaptable to a wider range of scientific content. Backtranslation, especially with domain-specific tagging and iterative approaches, offers another potential path for significant improvements. Additionally, advanced parallel sentence filtering techniques, like those offered by Bicleaner AI, can ensure high-quality training data. By incorporating these advancements, we can contribute to NMT models that effectively bridge the communication gap in the global scientific community.

Acknowledgements

This work was created within the Scilake¹¹ project. We are grateful to the Scilake project for providing the resources and support that made this work possible. This project has received funding from the European Union’s Horizon Europe framework programme under grant agreement No. 101058573.

References

- Altbach, P. G. 2007. The imperial tongue: English as the dominating academic language. *Economic and Political Weekly*, page 3608–3611.
- Artetxe, Mikel and Holger Schwenk. 2019a. Margin-based parallel corpus mining with multilingual sentence embeddings. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3197–3203.
- Artetxe, Mikel and Holger Schwenk. 2019b. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *Transactions of the association for computational linguistics*, 7:597–610.
- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Byrne, Jody. 2012. *Scientific and Technical Translation Explained*. Routledge.

¹¹<https://scilake.eu/>

- Edunov, Sergey, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. *arXiv preprint arXiv:1808.09381*.
- Fiorini, S, A Tezcan, T Vanallemeersch, S Szoc, K Migdisi, L Meeus, and L Macken. 2023. Translations and open science. In *Proceedings of the International Conference HiT-IT*, pages 41–51.
- Junczys-Dowmunt, Marcin, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, et al. 2018. Marian: Fast neural machine translation in c++. *arXiv preprint arXiv:1804.00344*.
- Ott, Myle, Sergey Edunov, David Grangier, and Michael Auli. 2018. Scaling neural machine translation. *arXiv preprint arXiv:1806.00187*.
- Papavassiliou, Vassilis, Sokratis Sofianopoulos, Prokopis Prokopidis, and Stelios Piperidis. 2018. The ilsp/arc submission to the wmt 2018 parallel corpus filtering shared task. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 928–933.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Popović, Maja. 2015. chrF: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.
- Popović, Maja. 2016. chrF deconstructed: beta parameters and n-gram weights. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 499–504.
- Popović, Maja. 2017. chrF++: words helping character n-grams. In *Proceedings of the second conference on machine translation*, pages 612–618.
- Post, Matt. 2018. A call for clarity in reporting bleu scores. *arXiv preprint arXiv:1804.08771*.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. *arXiv preprint arXiv:2009.09025*.
- Roussis, Dimitrios, Vassilis Papavassiliou, Prokopis Prokopidis, Stelios Piperidis, and Vassilis Katsouras. 2022. Scipar: A collection of parallel corpora from scientific abstracts. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 2652–2657.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2015. Improving neural machine translation models with monolingual data. *arXiv preprint arXiv:1511.06709*.
- Soares, Felipe, Gabrielli Harumi Yamashita, and Michel Jose Anzanello. 2018. A parallel corpus of theses and dissertations abstracts. In *Computational Processing of the Portuguese Language: 13th International Conference, PROPOR 2018, Canela, Brazil, September 24–26, 2018, Proceedings 13*, pages 345–352. Springer.
- Soares, Felipe, Viviane Pereira Moreira, and Karin Becker. 2019. A large parallel corpus of full-text scientific articles. *arXiv preprint arXiv:1905.01852*.
- Stockemer, D., Wigginton M.J. 2019. Publishing in english or another language: An inclusive study of scholar’s language publication preferences in the natural, social and interdisciplinary sciences. *Scientometrics*, 118:645–652.
- Tiedemann, Jörg and Santhosh Thottingal. 2020. OPUS-MT – building open translation services for the world. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 479–480, Lisboa, Portugal, November. European Association for Machine Translation.
- Tiedemann, Jörg. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC’12)*, pages 2214–2218, Istanbul, Turkey, May. European Language Resources Association (ELRA).
- Tiedemann, Jörg. 2020. The tatoeba translation challenge – realistic data sets for low resource and multilingual MT. In *Proceedings of the Fifth Conference on Machine Translation*, pages 1174–1182, Online, November. Association for Computational Linguistics.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30.

Appendix A

List of the 62 Processed Repositories

Name of Repository	URL
ESTUDO GERAL - Digital Repository of the University of Coimbra	https://estudogeral.uc.pt/
Repositorio da Universidade de Lisboa	https://repositorio.ul.pt/
RepositóriUM - Institutional Repository of the University of Minho	https://repositorium.sdum.uminho.pt/
Repositório Científico Lusófona	https://recil.ensinolusofona.pt/
Repositório da Universidade da Madeira	https://digituma.uma.pt/
RIA - Repositorio Institucional - Universidade de Aveiro	https://ria.ua.pt/
RUN - Repositorio Universidade Nova	https://run.unl.pt/
Repositório Institucional da Universidade Fernando Pessoa	https://bdigital.ufp.pt/
Repositorio Aberto - Universidade Aberta	https://repositorioaberto.uab.pt/
Repositorio de UTAD - Universidade de Tras-os-Montes e Alto Douro	https://repositorio.utad.pt/
Biblioteca Digital do IPB - Instituto Politecnico de Braganca	https://bibliotecadigital.ipb.pt/
Dadun - Open Access Institutional Repository of the University of Navarra	https://dadun.unav.edu/
UIBrepository - Universitat de les Illes Balears	https://dspace.uib.es/
Archivo Digital - Universidad Politécnica de Madrid	https://oa.upm.es/
Dipòsit Digital de la Universitat de Barcelona	http://diposit.ub.edu/
RODERIC - Repository of Universitat de Valencia	https://roderic.uv.es/
UCrea - Repositorio Abierto de la Universidad de Cantabria	https://repositorio.unican.es/
Gestión del Repositorio Documental de la Universidad de Salamanca (GREDOS)	https://gredos.usal.es/
Repository of Universitat Oberta de Catalunya	https://openaccess.uoc.edu/
RIULL - Repositorio Institucional Universidad de La Laguna	https://riull.ull.es/
idUS - Deposito de Investigacion Universidad de Sevilla	https://idus.us.es/
UVaDOC Repositorio Documental de la Universidad de Valladolid	https://uvadoc.uva.es/
UPF Digital Repository (Universitat Pompeu Fabra)	https://repositori.upf.edu/
RiuNet - Institutional Repository of the Politechnical University of Valencia - UPV	https://riunet.upv.es/
RUC - Repositorio Universidade de Coruna	https://ruc.udc.es/
UPCommons - Universitat Politècnica de Catalunya Barcelonatech	https://upcommons.upc.edu/
UAM - Repositorio institucional de la Universidad Autonoma de Madrid	https://repositorio.uam.es/
DUGiDocs - Institutional Repository of the Universitat de Girona	https://dugi-doc.udg.edu/
Repositori Universitat Jaume I	https://repositori.uji.es/
AccedaCRIS - Repository of Universidad de Las Palmas de Gran Canaria	https://accedacris.ulpgc.es/
DDD - Dipòsit Digital de Documents de la UAB (Universitat Autònoma de Barcelona)	https://ddd.uab.cat/
Tesis Doctorals en Xarxa	https://www.tdx.cat/
Repositorio Académico de la Universidad de Chile	https://repositorio.uchile.cl/
Biblioteca Institucional - Universidad Andres Bello	https://repositorio.unab.cl/
Repositorio Institucional UCA (Pontificia Universidad Católica Argentina)	https://repositorio.uca.edu.ar/
DSpace PUCP (Pontificia Universidad Católica del Perú)	https://tesis.pucp.edu.pe/
Repositorio Institucional - Universidad Nacional de Ingeniería	http://cybertesis.uni.edu.pe/
Repositorio Académico UPC (Universidad Peruana de Ciencias Aplicadas)	https://repositorioacademico.upc.edu.pe/
Repositorio Institucional Séneca - Universidad de los Andes	https://repositorio.uniandes.edu.co/
Institutional Repository of Universidad Nacional de Colombia	https://repositorio.unal.edu.co/
Institutional Repository - Pontificia Universidad Javeriana	https://repositorio.javeriana.edu.co/
Universidade de Brasília - Institutional Repository	https://repositorio.unb.br/
Universidade Federal de Santa Catarina - Institutional Repository	https://repositorio.ufsc.br/
Universidade Federal da Paraíba - Institutional Repository	https://repositorio.ufpb.br/
Repositório Institucional da Universidade Tecnológica Federal do Paraná (RIUT)	http://repositorio.utfpr.edu.br/
Repositório Institucional UNESP (Universidade Estadual Paulista)	https://repositorio.unesp.br/
LUME - Repositorio Digital UFRGS (Universidade Federal do Rio Grande do Sul)	https://lume.ufrgs.br/
Repositorio Institucional - Universidade Federal do Rio Grande do Norte (UFRN)	https://repositorio.ufrn.br/
Repositório Institucional UNIFESP (Universidade Federal de São Paulo)	https://repositorio.unifesp.br/
ATTENA - Repositório Digital da UFPE (Universidade Federal de Pernambuco)	https://repositorio.ufpe.br/
Repositorio Institucional da UFBA - Universidade Federal da Bahia	https://repositorio.ufba.br/
Institutional Repository PUCRS - Pontificia Universidade Católica do Rio Grande do Sul	https://repositorio.pucrs.br/
Repositório Institucional - Universidade Federal de Uberlândia	https://repositorio.ufu.br/
Locus Repositório Institucional da UFV (Universidade Federal de Viçosa)	https://www.locus.ufv.br/
DIAL - Research Publications of Université Catholique de Louvain	https://dial.uclouvain.be/pr/boreal/
Papyrus Institutional Repository - Université de Montreal	https://papyrus.bib.umontreal.ca/
Savoirs UdeS - Université de Sherbrooke	https://savoirs.usherbrooke.ca/
Toulouse Capitole Publications - Université Toulouse 1 Capitole	https://publications.ut-capitole.fr/
RED de Repositorios Latinoamericanos	https://repositorioslatinoamericanos.uchile.cl/
Portal de Revistas Academicas Chilenas	https://revistaschilenas.uchile.cl/
HAL Theses - Theses en Ligne	https://theses.hal.science/
HAL Open Science	https://hal.science/

Table 4: Processed Repositories Names and URLs

Appendix B

Example Configuration File for Processing the HTMLs

The configuration files used for each one of the repositories that we processed facilitate the structured extraction of metadata. The regex patterns specified within the configuration file are used to match the HTML structure of the academic records, while the minimum lengths specify the minimum characters that a title or abstract must have so as to be considered valid. Additionally, the file specifies the targeted languages for which we extract texts.

Below is the JSON configuration file for the DIAL UCLouvain repository:

```
{
  "abstracts_regex": ".*publication-metadata.*",
  "abstracts_min_len": 20,
  "titles_regex": ".*citation_title.*",
  "titles_min_len": 20,
  "keywords_regex": ".*Keywords.*",
  "authors_regex": ".*citation_author.*",
  "publishers_regex": ".*Affiliation.*|.*Publisher.*",
  "date_available_regex": ".*Publication date.*|.*Defense date.*",
  "journal_regex": ".*citation_journal_title.*|.*citation_dissertation_institution.*",
  "bibliographic_citation_regex": ".*Bibliographic reference.*",
  "document_language_regex": ".*Language.*",
  "link_html_regex": ".*Permanent URL.*",
  "link_pdf_regex": ".*citation_pdf_url.*",
  "document_type_regex": ".*Document type.*",
  "license_regex": ".*Access type.*",
  "URI_regex": ".*Permanent URL.*",
  "targeted_langs": [
    "en",
    "es",
    "pt",
    "fr"
  ]
}
```

Appendix C

Example JSON File Extracted from a Record

Our repository processing pipeling results into JSON files with the following structured metadata: abstracts, titles, repository name, links to the HTML and PDF versions of the documents, URI, license type, publication dates, document language and type (e.g., thesis), keywords, author names, publisher information, publication dates, journal titles, bibliographic citations, and the identified domains.

Below is an example JSON file extracted from a record, classified as belonging to the "Energy Research" domain, which originates from the Biblioteca Digital do IPB - Instituto Politécnico de Braganca repository:

```
{
  "abstracts": {
    "en": "In this research is intended to analyse the expansion of the economic sector related to the development ways of renewable energy and the economic and financial performance of companies operating in this field. [...]",
    "pt": "Esta investigação pretende analisar a expansão do setor económico relacionado com o desenvolvimento das energias renováveis e os desempenhos económico e financeiro das empresas que operam nesse setor. [...]"
  },
  "titles": {
    "en": "The development ways of renewable energy: the economic and financial performance of firms in this sector in Armenia and OECD countries"
  },
  "repository": "bibliotecadigital-ipb-pt",
  "html_id": 14638,
  "link_html": "https://bibliotecadigital.ipb.pt/handle/10198/14638",
  "link_pdf": "https://bibliotecadigital.ipb.pt/bitstream/10198/14638/1/Tarakhchyan_Siranush.pdf",
  "uri": "http://hdl.handle.net/10198/14638",
  "license_link": "http://creativecommons.org/licenses/by-nc/4.0/",
  "license": "openAccess",
  "date_available": "2017-11-20T15:08:42Z",
  "document_language": "en",
  "document_type": "masterThesis",
  "keywords": [
    "Renewable energy (RE)",
    "Financial data",
    "Financial ratios",
    "Market price",
    "Environment",
    "Domínio/Área Científica::Ciências Sociais::Economia e Gestão"
  ],
  "authors": [
    "Tarakhchyan, Siranush"
  ],
  "publishers": [],
  "bibliographic_citation": "",
  "journal": "",
  "domain_keyword_count": {
    "cancer": 0,
    "energy": 6,
    "transportation": 0,
    "neuroscience": 0
  },
  "domain": "energy"
}
```

Towards Tailored Recovery of Lexical Diversity in Literary Machine Translation

Esther Ploeger* Huiyuan Lai* Rik van Noord* Antonio Toral*

*Department of Computer Science, Aalborg University, Denmark

*CLCG, University of Groningen, The Netherlands

espl@cs.aau.dk {h.lai, r.i.k.van.noord, a.toral.ruiz}@rug.nl

Abstract

Machine translations are found to be lexically poorer than human translations. The loss of lexical diversity through MT poses an issue in the automatic translation of literature, where it matters not only *what* is written, but also *how* it is written. Current methods for increasing lexical diversity in MT are rigid. Yet, as we demonstrate, the degree of lexical diversity can vary considerably across different novels. Thus, rather than aiming for the rigid *increase* of lexical diversity, we reframe the task as *recovering* what is lost in the machine translation process. We propose a novel approach that consists of reranking translation candidates with a classifier that distinguishes between original and translated text. We evaluate our approach on 31 English-to-Dutch book translations, and find that, for certain books, our approach retrieves lexical diversity scores that are close to human translation.

1 Introduction

With the introduction of neural machine translation (NMT), the performance of high-resource automatic translation has improved substantially. Especially since the introduction of the Transformer architecture (Vaswani et al., 2017), state-of-the-art NMT systems have outperformed previous approaches considerably (Lakew et al., 2018), with some works even claiming human parity (Popel et al., 2020). However, these claims are based mostly

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

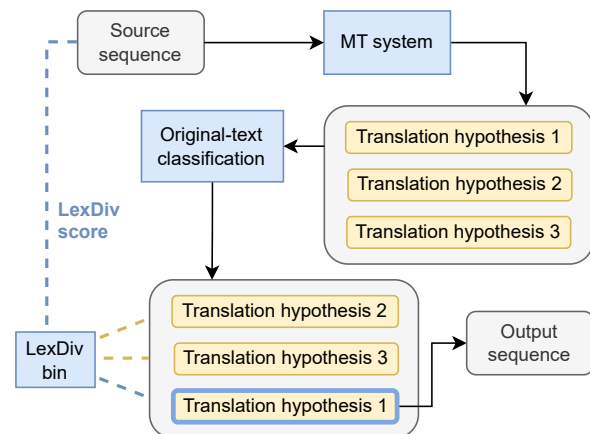


Figure 1: Reranking translation hypotheses based on the probability they are originally written in the target language, where the chosen rank is based on the lexical diversity score of the original book, and could be lower than the most lexically diverse option.

on accuracy and fluency measures, while style is often overlooked. In fact, according to expert evaluation, machine translation (MT) did actually not reach human parity (Toral et al., 2018; Fischer and Läubli, 2020). For instance, MT models have been found to exacerbate linguistic patterns that occur frequently, while underrepresenting patterns that are found less commonly (Vanmassenhove et al., 2019). As a result, automatically translated texts are found to be lexically poorer than human translations (HT). This ‘artificially impoverished language’ has previously been referred to as *machine translationese* (Vanmassenhove et al., 2021).

In this paper, we focus on the translation of novels. Contrary to technical domains, where meaning preservation is the main criterion for acceptable translations, literary translations have the additional criterion of style. This is because apart from meaning preservation (*what* is written), maintaining a certain reading experience (*how* it is writ-

ten) is vital for novels (Toral and Way, 2015). Importantly however, writing style (and linguistic complexity) can vary considerably between books. Some books contain repetitive language use, while others are written in embellished language (see Section 3). Current approaches that aim to mitigate the loss of lexical diversity do not accommodate this. State-of-the-art previous work (Freitag et al., 2019; Freitag et al., 2022) increases lexical diversity in a rigid way, not allowing for flexibility at inference time.

Contributions (i) We show that lexical diversity varies considerably across books, and argue that this should be taken into account in MT; (ii) We introduce a novel flexible method for recovering lexical diversity in MT, informed by the diversity of the original. (ii) We evaluate our method on 31 English novels which are translated to Dutch, and find that our approach is effective when it comes to book-tailored promotion of lexical diversity.

2 Related Work

Literary MT NMT has been argued to hold potential for literary texts, for instance in assisting professional translators or improving the immediate accessibility of untranslated foreign language books (Matusov, 2019). However, MT has been shown to decrease lexical diversity (Vanmassenhove et al., 2019; Vanmassenhove et al., 2021). This is an issue, because literary works can be viewed as a special domain in translation. Typically, literary translators are expected to preserve not only literal elements from the source, such as the plot, but also some sense of creative value (Riera, 2022). In other words, a goal of literary translation could be to recreate the ‘aesthetic intentions or effects’ that are possibly present in the source book (Delabastita, 2011). Such ‘aesthetic intentions’ can for instance be voice and metaphor, but also repetition (Wright, 2016). Repetitive use of language is commonly a conscious choice by the writer, and has a function, such as drawing attention or establishing a pattern (Boase-Beier, 2011). Given that lexical diversity can be an intentional writing choice, it should be apparent that an approach that aims at recovering lexical diversity in MT should not be boundless. Therefore, it is our aim to inform recovery with the degree of relative lexical diversity of the source text.

Machine Translationese Following recommendations from Jiménez-Crespo (2023), we will largely refrain from using the term *translationese* in the rest of this paper. However, it is important to note that previous work that aims to increase lexical diversity in MT has mostly been framed as part of ‘machine translationese’ reduction (Freitag et al., 2019; Freitag et al., 2022; Dutta Chowdhury et al., 2022; Jalota et al., 2023). Translations have been found to differ from original texts in a number of ways. For one, Baker (1993) argues that human translations into a language tend to be lexically simpler than text originally written in that language. Automatic classification approaches have been effective in detecting this difference (Baroni and Bernardini, 2005; Koppel and Ordan, 2011; Volansky et al., 2015; Rabinovich and Wintner, 2015; Pylypenko et al., 2021). More recently, work has investigated linguistic differences between MT and HT (van der Werff et al., 2022). Thus, it seems that modelling characteristics of original versus translated texts has a direct link to lexical diversity. Previous work (Freitag et al., 2022) leveraged these detectable differences in their approach to increase the naturalness of output translations. We take inspiration from their lexical diversity evaluation methods, and implement their method as a baseline.

Reranking Methods Reranking hypotheses in text generation originated before the age of neural paradigms (Shen et al., 2004; Collins and Koo, 2005). In essence, reranking entails re-ordering the set of candidate outputs according to some criterion, with the aim of providing a final output that adheres better to that criterion. Such methods have been applied for various tasks, such as summarization (Liu and Liu, 2021) and semantic parsing (Yin and Neubig, 2019). In machine translation, previous approaches include discriminative reranking (Lee et al., 2021) and reranking with energy-based models (Arcadinho et al., 2022).

3 Why Recover Rather Than Increase Lexical Diversity?

In this paper, we argue for tailored recovery of lexical diversity. In this section, we first discuss support for this idea from the field of literary studies. Then, we provide empirical evidence by applying lexical diversity metrics to our test set.

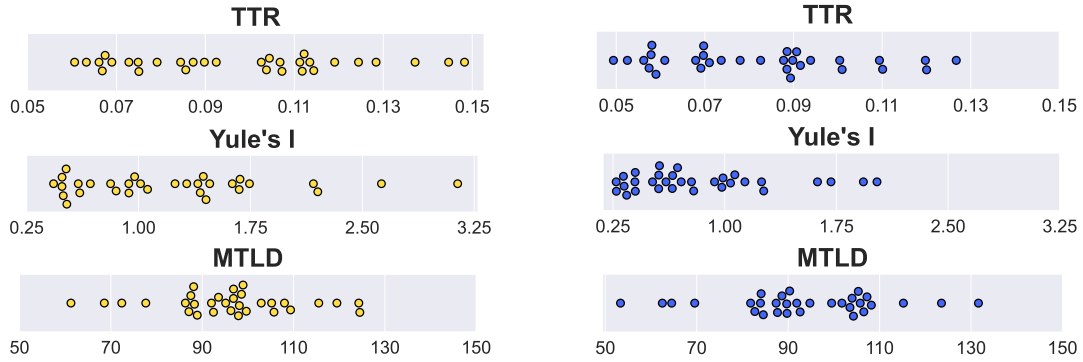


Figure 2: Range and spread of lexical diversity metrics for HT (left, yellow) and original English (right, blue).

3.1 Theoretical Support

Previous work on writing style in novels acknowledges that some books exhibit more lexical diversity than others. As an example, Heaton (1970) finds that no word in the original (i.e. English) version of *The Old man and the Sea* by Ernest Hemingway contains more than six syllables. Additionally, Hemingway tends to stick to particular words, even when there are more diverse options: in 184 situations of direct speech, he chooses to use the word ‘said’ 170 times instead of for example ‘asked’, ‘remarked’, ‘noticed’ or ‘yelled’. An example from the other end of the spectrum is James Joyce’s *Ulysses*. This work is known for its experimental techniques and unorthodox language use. Trotta (2014) illustrates this by highlighting Joyce’s use of neologisms, such as ‘He *smellsipped* the cordial juice’ and ‘Davy Byrne *smiledyawnednodded* all in one’. Moreover, Joyce repeatedly uses non-verbs as verbs, like in ‘I am *almosting* it.’ and even writes long sequences in unconventional spelling (*Ahbeesee defeegeee kelomen opeecue rustyouvee doubleyou. Boys are they?*). These examples make it clear that books can be written with vastly different ‘aesthetic intentions’. Thus, for preserving these intentions, MT approaches should not render them equally diverse in terms of lexicon.

3.2 Empirical Support

We empirically verify whether these findings hold for our data specifically, by estimating the lexical diversity of the 31 books in our test set, which we introduce in Section 5.1. We calculate three measures of lexical variety (type-token ratio; TTR, Yule’s I (Yule, 1944), and MTLD (McCarthy, 2005)) for each book in our test set. We further elaborate on these metrics in Section 6. Next, we apply the same metrics to the human translations

of those same books. Figure 2 shows that there is indeed a wide range of diversity across books, for both HT and original text. For example, in both settings, we find that the highest MTLD value is almost two times as large as the lowest. This emphasises why it is not our aim to generate the highest possible lexical diversity for every book. While we observe similar ranges and distributions in HT vs. original, the HT metrics are slightly higher. However, this does not necessarily mean that HT contains more embellished language. We note that the languages in our study, Dutch and English, are relatively similar (both in terms of genealogy and linguistic typology), but they differ in ways that can influence diversity metrics. For instance, Dutch contains compound nouns while English does not, making a higher TTR for Dutch more likely.

This discrepancy means that we cannot compare our Dutch MT to the original English book diversity directly. Instead, here we compare MT with HT. To verify whether this is sensible, we assess the relationship between HT and the English originals, by computing Pearson’s correlation on the corresponding diversity metrics. The results are listed in Table 1, and the corresponding regression plots are found in Appendix B. We observe strong correlations that are all statistically significant. This is important, because as the source diversity is a reliable indicator of HT diversity, it makes sense to use the source scores to approach HT (see Section 4).

Metric	Correlation coefficient	<i>p</i> -value
TTR	0.971	< 0.00001
Yule’s I	0.929	< 0.00001
MTLD	0.953	< 0.00001

Table 1: Pearson correlation coefficients for HT and OR lex-div metrics, rounded to three decimals.

4 Reranking Method

As illustrated in Figure 1, our approach consists of two parts: hypothesis generation and hypothesis reranking. Firstly, we generate the n best translation candidates for each source sentence in the test set with a vanilla domain-specific MT system (Section 5.1). Note that we decode all books separately, instead of concatenating all test set books. Then, for each book, we apply a classifier (Section 5.2) to the translation hypotheses and, through a softmax layer, obtain the probability for each candidate that it is an original Dutch sequence. Based on these probabilities, we rerank the translation candidates. In order to obtain the (expected) most lexically rich candidate, we would then choose the rank with the highest original-text probability. However, note that this simple approach is flexible in the sense that, instead of choosing the most original-like option, we have the option to choose a lower original-text rank.

We leverage this flexibility for tailoring rank selection to the lexical diversity of the original English book. First, for each original book, we calculate a *LexDiv* score, which consists of the average of the normalized TTR, Yule’s I and MTLTD scores (see Section 6). Then, we bin the books according to their *LexDiv* score, relative to the total distribution. That is, given a list that is sorted based on *LexDiv*, we categorize these into groups, where the number of groups depends on the number of n best candidates in decoding. For example, for $n = 5$, we bin the books into 5 different groups of 6 books (adding any remainders into the last bin). The bin per book corresponds to the original-text rank that is selected. As such, the selected rank for each book depends on the lexical diversity of its source, relative to the other books. Reranking translation candidates is a suitable solution to our task, because it accommodates flexibility, which is tunable at inference time. There is no need to train a separate model per diversity setting, saving computational expenses. Additionally, our approach is model-agnostic: reranking can be applied to any MT model that can generate multiple translation candidates.

5 Experimental set-up

5.1 Vanilla MT System

Data We use the dataset by Toral et al. (2024), which contains 531 books that were originally

written in English and manually translated into Dutch. We use 495 books for training, 5 for development and 31 as a test set. The genres of the books vary: they include literary fiction, popular fiction, non-fiction and children’s books from over 100 authors. We do not make a distinction between literary and ‘unliterary’ novels, as we believe this to be a subjective judgment.¹

Training Firstly, we align the sentences of the English and Dutch versions of each book using Vecalign (Thompson and Koehn, 2019). For the books in the test set, we manually discard sentences for which there existed no proper alignment, such as front matter sentences. Additionally, we discard sentences with a cosine distance higher than 0.7 (2.3% of all sentences). Then, we normalise all punctuation using the MOSES toolkit.² We then apply SentencePiece (Kudo and Richardson, 2018) subword segmentation to the data. For this, we train a SentencePiece unigram model with a joint vocabulary for both languages and a vocabulary size of 32,000.

We train a Transformer-based translation model using the Fairseq toolkit (Ott et al., 2019). More specifically, we use the *transformer_iwslt_de_en* architecture. This is a Transformer base model with 6 encoder and decoder layers and an embedding dimension of 512. During training, we use an Adam optimiser, a learning rate of $5e-4$, the loss function cross entropy with label smoothing 0.1 and the batch size is 64. Each model is trained until convergence with a patience of 3 epochs, using the BLEU score as a maximisation metric for finding the best checkpoint.

Decoding Strategies By default, we use beam search for decoding. Reranking approaches rely heavily on the diversity of the translation hypotheses: if the hypotheses are all very similar, reranking them is not likely to have a large effect. To ensure diverse hypotheses, we use a beam size of 20. Additionally, we experiment with decoding through diverse beam search (Vijayakumar et al., 2016). We follow Vijayakumar et al. (2016) by using 3 groups, with a beam size of 21. Beyond beam search, we investigate the effects of top-k and top-p sampling, with the default parameters and sampling size 10.

¹A full list of author names, titles, genres and publishing years of the test set books can be found in Appendix A, Table 8.

²<http://www.statmt.org/ Moses/>

System development (90%)				
Split	Orig.	# Books	# Sentences	# Words
Train (80%)	Dutch	1,291	8,576,756	10,425,656
	Other	1,291	12,470,149	165,263,466
Dev (10%)	Dutch	162	1,005,832	12,533,406
	Other	162	1,546,057	19,723,706
Test (10%)	Dutch	162	1,189,690	14,721,914
	Other	162	1,573,499	20,968,346
Original-text Classification (10%)				
Split	Orig.	# Books	# Sentences	# Words
Train (80%)	Dutch	143	982,114	11,528,789
	Other	143	139,0351	17,951,613
Test (20%)	Dutch	36	261,151	2,974,873
	Other	36	340,950	4,283,604
Total		3,588	29,336,549	376,130,733

Table 2: Monolingual data set division and size.

5.2 Original-Text Classification

Data We use a monolingual dataset of more than 7,000 Dutch books from varying original languages, authors and genres (Toral et al., 2024). For each book, we annotate whether it was originally written in Dutch.³ We discard 2,182 books for which the original language is unclear or that were not prose. We make sure to avoid overlap with the parallel data set by removing any books that are also part of the parallel data. Finally, we randomly sample 1,794 of the remaining 2,190 books as to match the total number of translated books, ensuring an equal distribution. In total, we are left with over 3,500 books and over 29M sentences. We further divide these into data for system development and data for original-text classification. We use this data for reproducing previous work (Freitag et al., 2022) and for training our classifier. Additionally, we translate the classifier section of the monolingual data set using a reverse-direction trained version of the vanilla MT system (NL → EN), and then perform round-trip-translation (RTT) back to Dutch with the vanilla MT system, to obtain an MT version of the monolingual classifier data. The full data size statistics and division in training, development and testing splits are listed in Table 2.

Training Currently, state-of-the-art performance for original-text detection is based on BERT (Devlin et al., 2019), as demonstrated by Pylypenko et al. (2021). We implement a similar system that distinguishes between original text and MT by train-

³The full annotation workflow can be found in Appendix C

ing a binary classification model. We fine-tune Dutch language model BERTje (de Vries et al., 2019). We train each model on the training split of the original-text classification data (see Table 2). We train models with batch size 128, accumulating gradients over 8 update steps, using the Adam optimiser (Kingma and Ba, 2015) with a learning rate of $3e-5$. We use early stopping (patience 3) if validation performance does not improve. On the held-out test set, the classifier achieves an accuracy of 85.9%. It obtains a precision of 90.6%, a recall of 80.2% and the F1 score is 85.0%.

5.3 Baselines

APE Freitag et al. (2019) introduced Automatic Post-Editing (APE) as a post-hoc method to increase the ‘naturalness’ of MT output. Following their approach, we train a post-processor that ‘translates’ synthetic Dutch sequences into more natural Dutch sequences. For training this system, we use the same data that was used to train the classifier (Section 5.2), consisting of RTT Dutch (which we use as source) and original Dutch (which we use as target). We train a model with the same architecture as the vanilla MT system. We apply the post-processor to the output of the vanilla MT system, in an attempt to obtain a translation with a lexical diversity that is closer to HT.

Tagging Our second baseline is based on Freitag et al. (2022). We train an MT system that learns to differentiate between original and translated text during training. This method requires both translated and original Dutch target samples. The translated target samples are found in our parallel dataset. We use the same original Dutch samples that are used in training the translationese classifier. Following Freitag et al. (2022), we then prepend $\langle orig \rangle$ to the English source sentences that have original Dutch on the target side, and $\langle trans \rangle$ for the source sentences that have translated Dutch. We train an MT system (same parameters as vanilla MT) on this data set. For inference, we prepend the source with $\langle orig \rangle$, which prompts the model to produce a translation that exhibits characteristics that are often found in original Dutch. Note that, in contrast to APE, this method cannot be applied post-hoc.⁴

⁴Note that our implementation differs from Freitag et al. (2022) in that they automatically differentiate natural and unnatural samples from a large parallel corpus using contrasting language models.

6 Evaluation

We introduce three classes of metrics. Firstly, we look at general text metrics, which are commonly used for evaluating lexical diversity. Secondly, we use translation-specific metrics. Lastly, we evaluate the general translation quality.

6.1 General Text Metrics

TTR The type-token ratio is the ratio of types (set of words) to tokens (actual words). A higher TTR indicates that more (different) words are used, which in turn indicates a higher lexical diversity. While this method is known to be influenced by the length of the text it is applied to, we report it because it is easy to interpret and widely used.

Yule’s I As a metric that is less sensitive to variation in text length, we use Yule’s I (Yule, 1944). We calculate this value as stated in Equation 1, where V is the size of the vocabulary (number of types) and $t(i, N)$ denotes the frequency of types which occur i times in a sample of length N .

$$\text{Yule's I} = \frac{V^2}{\sum_{i=1}^V i \times t(i, N) - V} \quad (1)$$

MTLD As an additional metric that has proven to be robust to document length variety, we use the measure of textual lexical diversity (MTLD), which is sequentially calculated as the ‘average length of sequential word strings in a text that maintain a given TTR value’ (McCarthy, 2005). We use the same TTR threshold (0.72) as Vanmassenhove et al. (2021).

We calculate these values using the *LexicalRichness* Python library (Shen, 2022).

6.2 Translation-specific Metrics

Vanmassenhove et al. (2021) introduce a novel automatic evaluation method for measuring lexical diversity in translations: Synonym Frequency Analysis (SFA). It provides an insight into the diversity of lexical choices in translations. For English words that have multiple translations in Dutch, it takes into account the frequency of these translation options. We re-implement this method, as it was not implemented for our language pair before. We first lemmatise each word in the source (English) side of our test set, using SpaCy (*nl_core_news_lg*).⁵ Next, we extract all possible translation options for the English adjectives,

⁵https://spacy.io/models/nl#nl_core_news_lg

nouns and verbs by using a English-to-Dutch bilingual dictionary.⁶ Next, for each translation option, we count the number of occurrences in the MT output for each system. The result is a vector which contains the occurrence frequency of each translation synonym for an English word.

PTF The primary translation frequency (PTF) is the average percentage (over all relevant source words) of times the most frequent translation option was chosen, from all translation options. The assumption is that if the output contains more secondary candidates, the text is more lexically diverse. We report the average PTF of all source words.

CDU The CDU is the cosine distance between the output vector for each source word and a vector of the same length with an equal distribution for each translation option (with the same total). We take the average CDU over all relevant source words to compute a final CDU.

SynTTR Lastly, we compute the SynTTR by dividing the number of types (the length of the set of all translation options) by the number of tokens (the sum of all translation options vectors).

6.2.1 Translation Quality

We also calculate a general measure of translation quality, because the ‘naturalness’ of a translation does not necessarily imply that a translation is a faithful representation of the source. A randomly generated string sequence might be very lexically diverse, but likely does not carry the source meaning. Firstly, we calculate BLEU (Papineni et al., 2002), as implemented in SacreBLEU (Post, 2018). We use the default settings, which are case-sensitive. Secondly, to account for the fact that BLEU does not necessarily evaluate meaning preservation, we additionally evaluate with COMET (Rei et al., 2020). English and Dutch are relatively high-resource languages, so we can use multilingual language embeddings. We report *comet-score*, calculated with the default *wmt22-comet-da*. Still, it should be noted that these automatic metrics do not necessarily correlate strongly with human judgements, especially for literary translation.

⁶We use the dictionary from <https://freedict.org/downloads>. As an example, for the English adjective *touching*, we find as Dutch translations: *ontroerend*, *aangrijpend*, *emotioneel*, *treffend*, *roerend* and *aandoenlijk*.

Approach	TTR \uparrow	Yule’s I \uparrow	MTLD \uparrow	PTF \downarrow	CDU \downarrow	SynTTR \uparrow	BLEU \uparrow	COMET \uparrow
HT	0.098	1.226	96.05	0.817	0.549	0.042	-	-
Vanilla MT	0.089	0.951	90.21	0.832	0.550	0.040	32.32	0.824
APE	0.092	0.985	90.59	0.827	0.554	0.041	30.39	0.808
Tagging	0.095	1.111	94.08	0.829	0.550	0.041	31.33	0.807
Tailored RR ($n=5$)	0.091	1.002	92.46	0.829	0.552	0.041	30.92	0.815
Tailored RR ($n=10$)	0.091	1.013	93.26	0.829	0.547	0.041	30.07	0.810
Tailored RR ($n=20$)	0.092	1.010	93.27	0.830	0.558	0.041	28.98	0.802
Tailored RR (<i>Top-k</i>)	0.101	1.286	104.25	0.815	0.559	0.043	21.21	0.745
Tailored RR (<i>Top-p</i>)	0.092	1.017	91.21	0.828	0.552	0.041	29.97	0.808
Tailored RR (<i>DBS</i>)	0.092	1.010	92.70	0.828	0.553	0.040	29.36	0.805

Table 3: Scores averaged across books, where RR stands for reranking. We provide results for multiple decoding strategies. Beam size is 20. Scores closest to HT are in bold font.

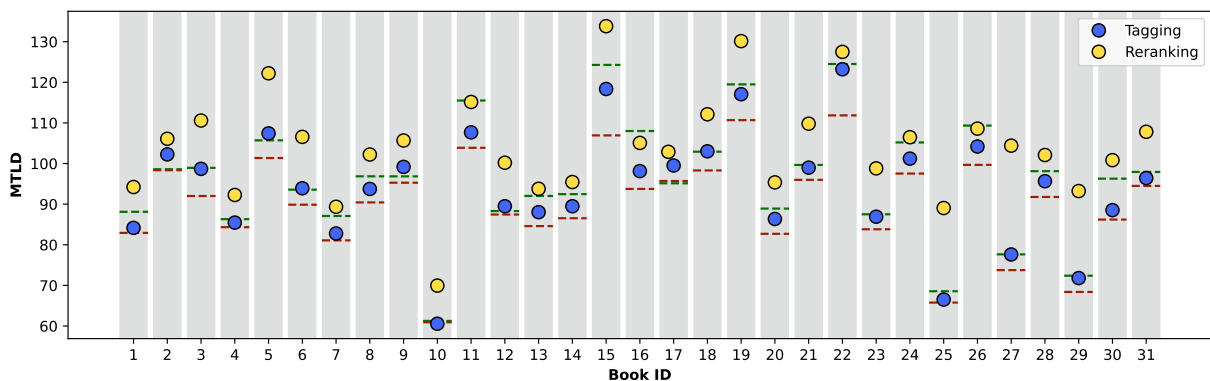


Figure 3: Per-book comparison of MTLD between the (rigid) tagging baseline and (tailored) reranking method, where green dotted lines are HT scores, and red dotted lines represent vanilla MT.

7 Results and Analysis

7.1 Quantitative Results

We first discuss the results over all books. Table 3 shows the average results of measuring lexical diversity and general translation quality across the various approaches. We find that vanilla MT indeed produces lexically poorer translations than HT, according to all our metrics. While the scores of the APE baseline remain close to vanilla MT, our tailored reranking approach retrieves a lexical diversity that is closer to HT. This suggests that our method is a suitable alternative for post-hoc editing, given that one has access to the MT model for generating translation hypotheses. The tagging baseline, which cannot be applied post-hoc, retrieves an MTLD and CDU that is on average closest to HT. Importantly though, it should be noted that reranking and tagging are not mutually exclusive: one could apply reranking to the tagging baseline to increase or decrease lexical diversity further, where desired. When we compare decoding strategies of the tailored reranking method,

we first observe that using diverse beams search and choosing a larger n retrieves at most slightly more diversity. Especially top-k decoding retrieves a much higher lexical diversity. However, tailored reranking comes with a compromise in terms of translation quality metrics.

Next, we demonstrate that these averages omit a more fine-grained view. Figure 3 shows the difference in MTLD per book between vanilla MT, HT, tagging and our most diverse reranking system, based on top-k sampling, which is tailored to the LexDiv score of the original English book.⁷ Our method renders almost every single book more lexically diverse than the tagging baseline. In some cases, this makes the results closer to HT in terms of lexical diversity (e.g. 7, 13, 14, 16). However, especially in cases where vanilla MT and HT are close already, this is not always true (e.g. 1, 3, 5).

⁷A similar figure with the posthoc baseline APE instead of tagging is shown in Appendix D.

Ex. #	Approach	Text
1	Source	The kid had no mother.
	HT	Dat joch heeft geen moeder gehad.
	Vanilla MT	Het kind had geen moeder.
	Tagging	De jongen had geen moeder.
	Tailored RR	Het joch had geen moeder.
2	Source	He shipped his oars and brought a small line from under the bow.
	HT	Hij haalde de riemen in en pakte een kleine lijn die voor in de boot lag.
	Vanilla MT	Hij trok zijn riemen aan en haalde een klein lijntje onder de boeg vandaan.
	Tagging	Hij verscheurde zijn riemen en haalde een klein streepje onder de boeg vandaan.
	Tailored RR	Hij haalde zijn riemen en trok er een kleine lijn voor onder de boot vandaan.
3	Source	In long shaky strokes Sargent copied the data.
	HT	In lange beverige halen kopieerde Sargent de gegevens.
	Vanilla MT	Met lange, bevende slagen kopieerde Sargent de gegevens.
	Tagging	Met lange bevende halen kopieerde Sargent de gegevens.
	Tailored RR	Met lange beverige halen schreef Sargent de data over .

Table 4: Examples to highlight surface-level differences between the systems’ output translations, where Tailored RR uses top-k sampling.

7.2 Surface-level Inspection

The output translations were inspected by a native speaker. Table 4 shows three examples of how translations differ between vanilla MT, tagging and tailored reranking (with top-k sampling). In Example 1 (from book 1, *Sunset Park*), we see that the English noun ‘kid’ is translated as *joch* (‘boy’) in the human translation, which is less common than the vanilla MT’s *kind* (‘child’) and tagging’s *jongen* (‘boy’). This is recovered by our tailored reranking system, which uses *joch* too.

Example 2 is taken from book 10, *The Old Man and the Sea*, which has low lexical diversity by default (see Section 3). This is not taken into account by the tagging baseline: the English ‘shipped’ is translated as a less common (and wrong) *verscheurde* (‘shredded’). The tailored reranking system (*haalde*, ‘brought’) is closest to HT (*haalde in*, ‘brought in’). Additionally, the tagging baseline wrongly translates the English ‘line’ as *streepje* (‘small stripe’), while tailored reranking (*lijn*, ‘line’) is again identical to HT. This case illustrates that choosing a more common translation synonym, which may for instance results in a lower PTF, may for some books be closer to HT.

By contrast, in Example 3 from the more lexically diverse *Ulysses* (book 15), the tagging baseline stays closer to vanilla MT: both translate

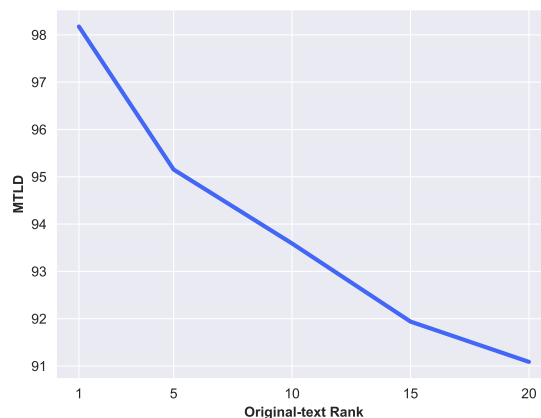


Figure 4: Change in MTLD for choosing different ranks, where beam size is 20 and $n = 20$.

‘shaky’ as *bevend* (‘trembling’). Tailored reranking outputs *beverig* (‘shaky’), which is again recovering the HT. Furthermore, tailored reranking deviates from all other systems (and HT) by translating ‘copied’ into the translation synonym *schreef over* (copying something by writing). This case may illustrate why the tailored reranking based on top-k sampling surpasses the other systems in the overall metrics.

7.3 Ranks and Lexical Diversity

So far, we have assumed that reranking based on the probability of a candidate being original text leads to more lexically diverse output translations.

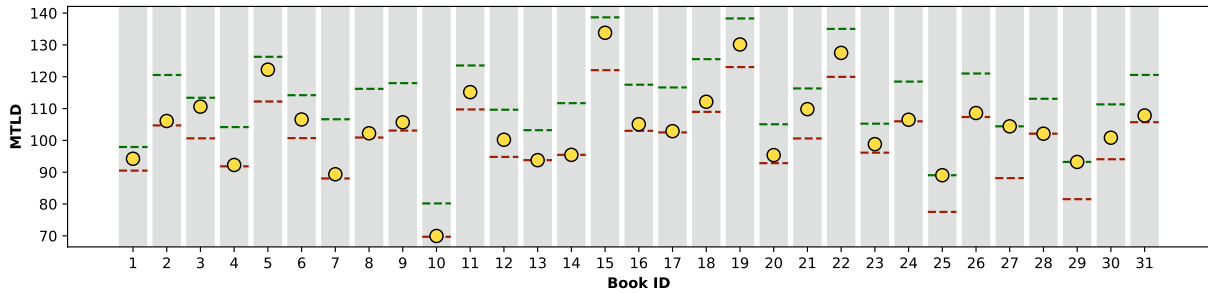


Figure 5: MTLD for highest (green), lowest (red) and tailored (yellow) original-text rank.

Here, we verify whether choosing a lower probability of a candidate being original, actually implies lexically poorer output translations (Figure 4). For the vanilla MT system with beam size 20 and $n = 20$, we first calculate the original-text probability for each translation hypothesis. Similar to reranking, we sort the hypotheses according to this probability. Then, instead of binning, we choose the n^{th} rank, and calculate lexical diversity of the output. Figure 4 shows the change in MTLD scores for choosing a lower diversity rank. We observe that indeed, choosing a lower rank retrieves lower diversity (note that there, a higher rank represents a smaller original-text probability). This trend holds for TTR and Yule’s I as well (see Appendix E).

7.4 Tailoring and Lexical Diversity

To further demonstrate the effect of a *tailored* approach in lexical diversity, we compare MTLD scores of a top-k reranking system that always outputs the highest original-text probability, with the same system that always outputs the lowest, and a tailored version. Figure 5 shows the results. Firstly, we observe that, in every case, choosing a rank that represents lower original-text probability retrieves a lower MTLD score than choosing the opposite. This corroborates the findings from the previous section. Next, we look into how the tailored reranking affects the output lexical diversity. In Section 3, we used *The Old Man and the Sea* (book 10) as an example of a book with a low default lexical diversity. We observe that our tailored reranking system outputs the lowest original-text probability rank for this book, resulting in a lower MTLD score. For the example from Section 3 of a lexically rich book, *Ulysses* (book 15), our tailored system outputs a rank with a original-text probability higher than the minimum, thus retrieving an MTLD score that is higher. This shows that tailoring is at least somewhat intuitive.

8 Conclusion

We have argued for flexible recovery of lexical diversity in literary MT. We showed that default diversity varies per book in our dataset, and that this lexical diversity is partially lost through MT. We presented the first approach towards tailored rescoring of translation candidates, which matches HT more closely than previous baselines for some books. Future work could explore how our method can be combined with previous work, as it is in principle model-agnostic. Investigations with document-level translation, instead of sentence-level translation only, could provide additional insights. Furthermore, it may be useful to address this task at an even finer-grained level, by exploring diversity reranking on a sequence-level, instead of a book-level.

Limitations

In this paper, we addressed the increase of lexical diversity in literary MT. However, it should be noted that this does not encompass writing style as a whole. We evaluated our approach on one high-resource language pair that consist of relatively similar languages, in one translation direction. For the domain of literary translation, we find this to be difficult to avoid. Still, experiments with more languages and resource-scenarios may retrieve interesting results. Moreover, while our data is transparent in the sense that we know and can explain exactly what it contains, we cannot distribute the data ourselves because of copyright. Lastly, we acknowledge that large-scale human evaluation could give useful insights into the differences between the systems.

Acknowledgements

This work was supported by a *Semper Ardens: Accelerate research grant (CF21-0454)* from the

Carlsberg Foundation. We thank the Center for Information Technology of the University of Groningen for their support and for providing access to the Peregrine an Hábrók high performance computing cluster.

References

- Arcadinho, Samuel David, David Aparicio, Hugo Veiga, and Antonio Alegria. 2022. T5QL: Taming language models for SQL generation. In Bosselut, Antoine, Khyathi Chandu, Kaustubh Dhole, Varun Gangal, Sebastian Gehrmann, Yacine Jernite, Jekaterina Novikova, and Laura Perez-Beltrachini, editors, *Proceedings of the 2nd Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 276–286, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Baker, Mona. 1993. Corpus linguistics and translation studies—implications and applications. In *Text and Technology*, page 233. John Benjamins.
- Baroni, Marco and Silvia Bernardini. 2005. A New Approach to the Study of Translationese: Machine-learning the Difference between Original and Translated Text. *Literary and Linguistic Computing*, 21(3):259–274, 08.
- Boase-Beier, Jean. 2011. *A critical introduction to translation studies*. Bloomsbury Publishing.
- Collins, Michael and Terry Koo. 2005. Discriminative reranking for natural language parsing. *Computational Linguistics*, 31(1):25–70.
- de Vries, Wietse, Andreas van Cranenburgh, Arianna Bisazza, Tommaso Caselli, Gertjan van Noord, and Malvina Nissim. 2019. BERTje: A Dutch BERT Model. arXiv:1912.09582, December.
- Delabastita, Dirk. 2011. Literary translation. *Handbook of translation studies*, 2:69–78.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Dutta Chowdhury, Koel, Richa Jalota, Cristina España-Bonet, and Josef Genabith. 2022. Towards debiasing translation artifacts. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3983–3991, Seattle, United States, July. Association for Computational Linguistics.
- Fischer, Lukas and Samuel Lübli. 2020. What’s the difference between professional human and machine translation? a blind multi-language study on domain-specific MT. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 215–224, Lisboa, Portugal, November. European Association for Machine Translation.
- Freitag, Markus, Isaac Caswell, and Scott Roy. 2019. APE at scale and its implications on MT evaluation biases. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 34–44, Florence, Italy, August. Association for Computational Linguistics.
- Freitag, Markus, David Vilar, David Grangier, Colin Cherry, and George Foster. 2022. A natural diet: Towards improving naturalness of machine translation output. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3340–3353, Dublin, Ireland, May. Association for Computational Linguistics.
- Heaton, CP. 1970. Style in the old man and the sea. *Style*, pages 11–27.
- Jalota, Richa, Koel Chowdhury, Cristina España-Bonet, and Josef van Genabith. 2023. Translating away translationese without parallel data. In Bouamor, Houda, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7086–7100, Singapore, December. Association for Computational Linguistics.
- Jimenez-Crespo, Miguel A. 2023. “translationese” (and “post-editeese”?) no more: on importing fuzzy conceptual tools from translation studies in MT research. In Nurminen, Mary, Judith Brenner, Maarit Koponen, Sirkku Latomaa, Mikhail Mikhailov, Frederike Schierl, Tharindu Ranasinghe, Eva Vanmassenhove, Sergi Alvarez Vidal, Nora Aranberri, Mara Nunziatini, Carla Parra Escartín, Mikel Forcada, Maja Popovic, Carolina Scarton, and Helena Moniz, editors, *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 261–268, Tampere, Finland, June. European Association for Machine Translation.
- Kingma, Diederik P. and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In Bengio, Yoshua and Yann LeCun, editors, *3rd International Conference on Learning Representations (ICLR 2015)*.
- Koppel, Moshe and Noam Ordan. 2011. Translationese and its dialects. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 1318–1326, Portland, Oregon, USA, June. Association for Computational Linguistics.
- Kudo, Taku and John Richardson. 2018. Sentence-Piece: A simple and language independent subword

- tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium, November. Association for Computational Linguistics.
- Lakew, Surafel Melaku, Mauro Cettolo, and Marcello Federico. 2018. A comparison of transformer and recurrent neural networks on multilingual neural machine translation. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 641–652, Santa Fe, New Mexico, USA, August. Association for Computational Linguistics.
- Lee, Ann, Michael Auli, and Marc’Aurelio Ranzato. 2021. Discriminative reranking for neural machine translation. In Zong, Chengqing, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7250–7264, Online, August. Association for Computational Linguistics.
- Liu, Yixin and Pengfei Liu. 2021. SimCLS: A simple framework for contrastive learning of abstractive summarization. In Zong, Chengqing, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 1065–1072, Online, August. Association for Computational Linguistics.
- Matusov, Evgeny. 2019. The challenges of using neural machine translation for literature. In Hadley, James, Maja Popović, Haithem Afli, and Andy Way, editors, *Proceedings of the Qualities of Literary Machine Translation*, pages 10–19, Dublin, Ireland, August. European Association for Machine Translation.
- McCarthy, Philip M. 2005. *An assessment of the range and usefulness of lexical diversity measures and the potential of the measure of textual, lexical diversity (MTLD)*. Ph.D. thesis, The University of Memphis.
- Ott, Myle, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Popel, Martin, Marketa Tomkova, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar, and Zdeněk Žabokrtský. 2020. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. *Nature communications*, 11(1):1–15.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October. Association for Computational Linguistics.
- Pylypenko, Daria, Kwabena Amponsah-Kaakyire, Koel Dutta Chowdhury, Josef van Genabith, and Cristina España-Bonet. 2021. Comparing feature-engineering and feature-learning approaches for multilingual translationese classification. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8596–8611, Online and Punta Cana, Dominican Republic, November. Association for Computational Linguistics.
- Rabinovich, Ella and Shuly Wintner. 2015. Unsupervised identification of translationese. *Transactions of the Association for Computational Linguistics*, 3:419–432.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.
- Riera, Jorge Braga. 2022. Literatura-traducción. *Enciclopedia de Traducción e Interpretación*.
- Shen, Libin, Anoop Sarkar, and Franz Josef Och. 2004. Discriminative reranking for machine translation. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 177–184.
- Shen, Lucas. 2022. LexicalRichness: A small module to compute textual lexical richness.
- Thompson, Brian and Philipp Koehn. 2019. Vecalign: Improved sentence alignment in linear time and space. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1342–1348, Hong Kong, China, November. Association for Computational Linguistics.
- Toral, Antonio and Andy Way. 2015. Machine-assisted translation of literary text: A case study. *Translation Spaces*, 4(2):240–267.
- Toral, Antonio, Sheila Castilho, Ke Hu, and Andy Way. 2018. Attaining the unattainable? Reassessing

- claims of human parity in neural machine translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 113–123, Brussels, Belgium, October. Association for Computational Linguistics.
- Toral, Antonio, Andreas van Cranenburgh, and Tia Nijssen. 2024. *Literary-adapted machine translation in a well-resourced language pair: Explorations with More Data and Wider Contexts*, pages 27–52. Routledge.
- Trotta, Joe. 2014. Creativity, playfulness and linguistic carnivalization in James Joyce's *Ulysses*.
- van der Werff, Tobias, Rik van Noord, and Antonio Toral. 2022. Automatic discrimination of human and neural machine translation: A study with multiple pre-trained models and longer context. In *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, pages 161–170, Ghent, Belgium, June. European Association for Machine Translation.
- Vanmassenhove, Eva, Dimitar Shterionov, and Andy Way. 2019. Lost in translation: Loss and decay of linguistic richness in machine translation. In *Proceedings of Machine Translation Summit XVII: Research Track*, pages 222–232.
- Vanmassenhove, Eva, Dimitar Shterionov, and Matthew Gwilliam. 2021. Machine translationese: Effects of algorithmic bias on linguistic complexity in machine translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2203–2213, Online, April. Association for Computational Linguistics.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Vijayakumar, Ashwin K, Michael Cogswell, Ramprasad R Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2016. Diverse beam search: Decoding diverse solutions from neural sequence models. *arXiv preprint arXiv:1610.02424*.
- Volansky, Vered, Noam Ordan, and Shuly Wintner. 2015. On the features of translationese. *Digital Scholarship in the Humanities*, 30(1):98–118.
- Wright, Chantal. 2016. *Literary translation*. Routledge.
- Yin, Pengcheng and Graham Neubig. 2019. Reranking for neural semantic parsing. In Korhonen, Anna, David Traum, and Lluís Màrquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4553–4559, Florence, Italy, July. Association for Computational Linguistics.
- Yule, C Udney. 1944. *The statistical study of literary vocabulary*. Cambridge University Press.

A Test set novels

ID	Author	Title	Year Published	Genre
1	Paul Auster	Sunset Park	2010	Literary fiction
2	David Baldacci	Divine Justice	2008	Thriller, suspense
3	Julian Barnes	The Sense of an Ending	2011	Literary fiction
4	John Boyne	The Boy in the Striped Pyjamas	2006	Historical fiction
5	John le Carré	Our Kind of Traitor	2010	Thriller, spy fiction
6	Jonathan Franzen	The Corrections	2001	Literary fiction
7	Nicci French	Blue Monday: A Frieda Klein Mystery	2011	Thriller, suspense
8	William Golding	Lord of the Flies	1954	Literary fiction
9	John Grisham	The Confession	2010	Thriller, suspense
10	Ernest Hemingway	The Old Man and the Sea	1952	Literary fiction
11	Patricia Highsmith	Ripley Under Water	1991	Thriller, suspense
12	Khaled Hosseini	A Thousand Splendid Suns	2007	Literary fiction
13	John Irving	Last Night in Twisted River	2009	Literary fiction
14	E.L. James	Fifty Shades of Grey	2011	Erotic thriller
15	James Joyce	Ulysses	1922	Literary fiction
16	Jack Kerouac	On the Road	1957	Literary fiction
17	Stephen King	11/22/63	2011	Science-fiction
18	Sophie Kinsella	Shopaholic and Baby	2007	Popular literature
19	David Mitchell	The Thousand Autumns of Jacob de Zoet	2010	Historical fiction
20	George Orwell	1984	1949	Literary fiction
21	James Patterson	The Quickie	2007	Thriller, suspense
22	Thomas Pynchon	Gravity's Rainbow	1973	Historical fiction
23	Philip Roth	The Plot Against America	2004	Political fiction
24	J.K. Rowling	Harry Potter and the Deathly Hallows	2007	Fantasy
25	J.D. Salinger	The Catcher in the Rye	1951	Literary fiction
26	Karin Slaughter	Fractured	2008	Thriller, suspense
27	John Steinbeck	The Grapes of Wrath	1939	Literary fiction
28	J.R.R Tolkien	The Return of the King	1955	Fantasy
29	Mark Twain	Adventures of Huckleberry Finn	1884	Literary fiction
30	Oscar Wilde	The Picture of Dorian Gray	1890	Literary fiction
31	Irvin D. Yalom	The Spinoza Problem	2012	Historical fiction

Table 5: Information on test set books.

B Regression plots for human translation vs. original text lexical diversity

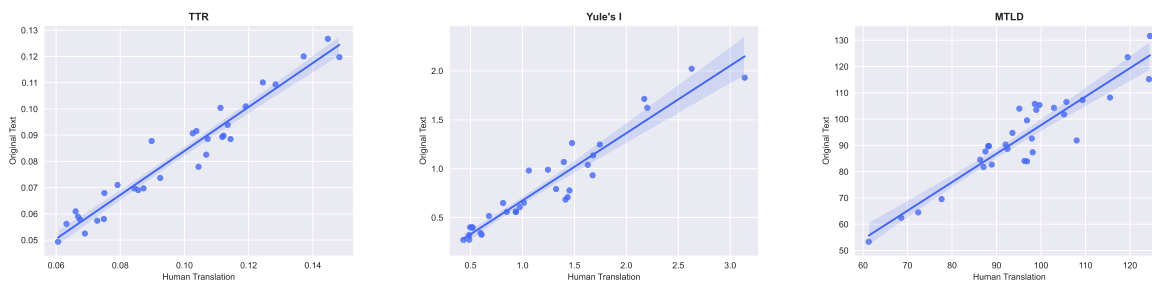


Figure 6: Regression plots for TTR, Yule's I and MTLD, with on the y-axis the scores for the original (English) versions, and on the x-axis those for human translations.

C Annotation workflow for monolingual Dutch books

1. Check whether the book is prose: we generally discard other forms of literature such as poetry and plays and annotate this in category 3 (no label).
2. Check whether the original language of the book is listed on the website of the National Dutch Library.⁸ If this is not the case:
 - (a) Check whether the language of the book is listed on the website of a Dutch reading community website.⁹
 - (b) If step (a) is also not conclusive: check whether more information on the author is available, for instance on a personal website where we can find the original titles.
 - (c) In case there is no reliable information available on the original language of a book, we discard the book (category 3: no label)
3. Book titles with Dutch as their original language are annotated with the label ‘1’ (category 1). Books that were written in a language other than Dutch were annotated with the label ‘0’ (category 2).

Special cases An interesting annotation case regards books from bilingual authors who learned Dutch at a later age, such as Kader Abdolah. In our current guidelines, we do not take this into account specifically; if originally written in Dutch, these books are annotated with category 1. We note that books that were translated to Dutch were not all originally written in English: other source languages in the data set include German, French and Spanish.

D Book-level MTLD comparison of APE and tailored reranking (top-k sampling)

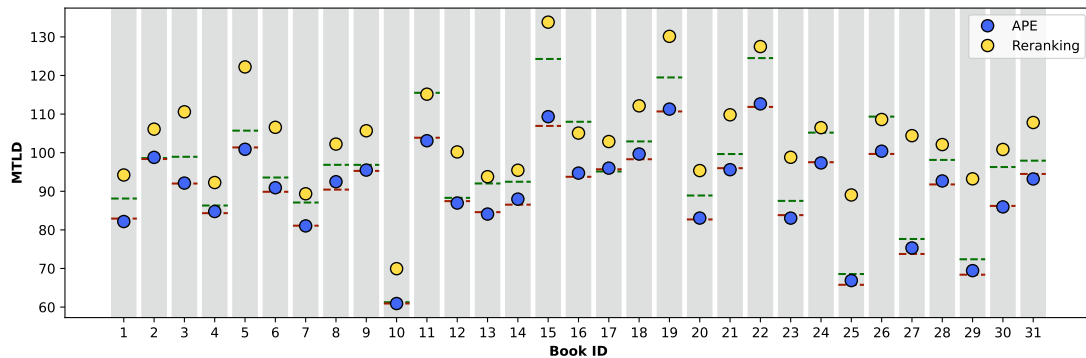


Figure 7: MTLD scores for APE and tailored reranking with top-k sampling, with on the y-axis the MTLD score for each book in our test set (x-axis).

E Lexical diversity according to ranks

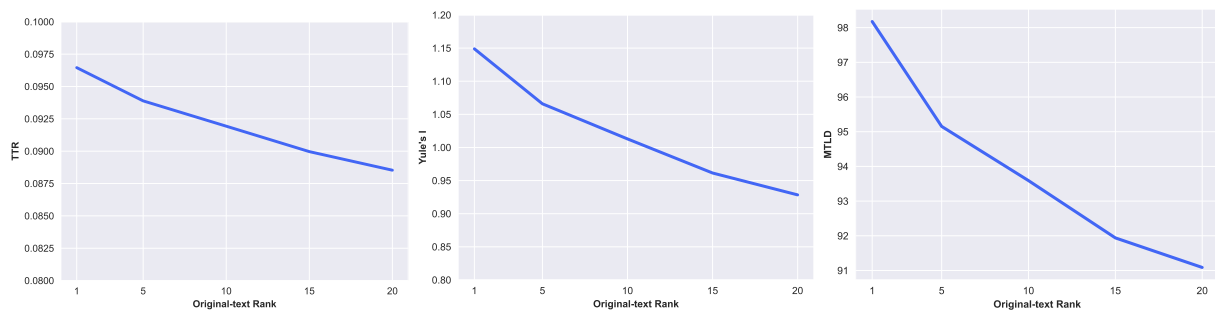


Figure 8: TTR, MTLD and Yule’s I according to original-text rank, where a higher rank represents smaller original-text probability.

⁸<https://www.bibliotheek.nl/>

⁹<https://www.hebban.nl/>

Enhancing Gender-Inclusive Machine Translation with Neomorphemes and Large Language Models

Andrea Piergentili,^{1,2} Beatrice Savoldi,² Matteo Negri,² Luisa Bentivogli²

¹University of Trento

²Fondazione Bruno Kessler

{apiergentili, bsavoldi, negri, bentivo}@fbk.eu

Abstract

Machine translation (MT) models are known to suffer from gender bias, especially when translating into languages with extensive gendered morphology. Accordingly, they still fall short in using gender-inclusive language, also representative of non-binary identities. In this paper, we look at gender-inclusive neomorphemes, neologistic¹ elements that avoid binary gender markings as an approach towards fairer MT. In this direction, we explore prompting techniques with large language models (LLMs) to translate from English into Italian using neomorphemes. So far, this area has been under-explored due to its novelty and the lack of publicly available evaluation resources. We fill this gap by releasing NEO-GATE,² a resource designed to evaluate gender-inclusive en→it translation with neomorphemes. With NEO-GATE, we assess four LLMs of different families and sizes and different prompt formats, identifying strengths and weaknesses of each on this novel task for MT.

1 Introduction

Machine translation (MT) has been found to be susceptible to gender bias, i.e. the tendency to produce default masculine outputs or stereotypical gender associations (Saunders et al., 2020; Savoldi

et al., 2021; Piazzolla et al., 2023) when gender information about human referents is absent. This is especially relevant when translating from notional gender languages like English, which express gender only through a limited set of elements (e.g., *he/she* pronouns), into grammatical gender target languages, such as Italian, German, and Spanish, which mark gender extensively in their morphology (e.g., en: “*My friends are rich*” → it: “*I miei amici sono ricchi*” [M] vs “*Le mie amiche sono ricche*” [F]). The consequences of this behavior are systematically harmful (Blodgett et al., 2020) to women, who risk being under-represented and stereotypically defined, and non-binary individuals, who are erased from representation or misgendered within binary gender linguistic frameworks (Misieki, 2020; Dev et al., 2021).

In light of this, in this paper we look at neologistic solutions – which are emerging from grassroots efforts to make language more inclusive – as a path towards gender-inclusive MT. Linguistic innovations such as neopronouns (e.g., en *ze* instead of *he/she*, sw *hen* instead of *han/hon*) and neomorphemes (e.g., it *-ə/-3*, es *-el/-es* in place of gendered inflectional morphemes) add new elements to the grammar and morphology to allow for the expression of non-binary gender identities or to convey gender neutrality (Bradley et al., 2019). To date, the use of neologistic solutions is not systematized yet, with alternative paradigms coexisting and new ones continuously emerging. The choice and use of one paradigm of neologistic devices (e.g., the neopronouns *xe/xem/xyr/xyrs/xemself* vs *ze/zir/zir/zirs/zirself*, etc.) depends on individuals’ identity and preferences in gender expression.

The use of neologistic devices in MT is still a largely unexplored research area, due to the novelty of this approach and to the lack of dedicated

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹Following (Rose et al., 2023), we refer to them as *neologistic* because of their linguistically innovative nature.

²Available at <https://huggingface.co/datasets/FBK-MT/Neo-GATE> under the CC-BY-4.0 licence.

EN I like being surrounded by my friends.

- A** Mi piace essere circondato dai miei amici.
B Mi piace avere persone amiche intorno a me.
[en: I like to have people who are friends around me]
C Mi piace che intorno a me siano presenti persone che considero mie amiche.
[en: I like that there are people around me who I consider my friends]
D Mi piace essere circondat* da* mie* amic*.
E Mi piace essere circondat_o da₃ mie₃ amic₃.

Table 1: Examples of en→it translations with no gender information in the source. Example A uses generic masculine formulations to refer to human beings (in bold), while the rest employ different gender-inclusive strategies (underlined). B and C use periphrases of different verbosity, while D and E employ different neomorpheme paradigms.

resources, which in turn is complicated by the unfixed nature of these solutions. Ideally, gender-inclusive MT research should factor in the multiplicity of paradigms that make up the landscape of neologistic devices (Lauscher et al., 2022). However, the unavailability of evaluation and training resources which can be adjusted to any paradigm is a bottleneck for the investigation of gender-inclusive MT. Also, neural MT has been proven to fail in handling neologistic gender-inclusive language (Lauscher et al., 2023). Looking at other options, LLMs’ ability to adapt to unseen tasks through in-context learning (Brown et al., 2020; Min et al., 2022) offers a viable path toward gender-inclusive MT without the need for extensive training data. Thus, in this work we investigate multilingual LLMs’ ability to adapt to new, inclusive morphological paradigms in translation.

To this aim, we *i*) release NEO-GATE, a benchmark to evaluate gender-inclusive en→it translation with any of the ever-emerging neomorpheme paradigms; *ii*) explore different prompting strategies for neologistic gender-inclusive MT, across three open and one commercial LLMs, and the two most popular Italian neomorpheme paradigms.

2 Background

Following evolving social and linguistic phenomena (Sendén et al., 2021; Waldendorf, 2023), there has been a rising demand for the integration of gender-inclusive language in natural language processing (NLP) technologies to make them inclusive of all gender identities (Dev et al., 2021). In MT, gender-neutral translation (GNT) was recently proposed as a gender-fair approach to make translation technologies less biased and more inclusive (Piergentili et al., 2023a). GNT consists in using gender-neutralization strategies, such as

epicene formulations, (e.g., ‘persone’ – en: *people* – in examples B and C in Table 1), to avoid expressing the gender of human beings in the target language. However, this approach has considerable limitations: *i*) it can result in verbose phrasings, as in example C, which are only acceptable in specific contexts and textual domains, namely formal and institutional communication (Piergentili et al., 2023b); *ii*) it is arguably impossible to translate some terms by applying these strategies in grammatical gender languages like Italian (e.g., kinship terms, such as *parent* → it *genitore/genitrice*) (Motschenbacher, 2014). Moreover, the use of circumlocutory language to avoid expressing gender is regarded as a form of *indirect* non-binary language (Attig and López, 2020), in that it conceals gender, while other, *direct* solutions emphasize it.

Indeed, innovative alternatives have been proposed by queer communities as well. Neologistic elements, such as neopronouns and neomorphemes, have emerged in notional gender languages, such as Swedish (Gustafsson Sendén et al., 2015) and English (McGaughey, 2020), as well as in grammatical gender languages, such as Spanish (Sarmiento, 2015), French (Kaplan, 2022), and German (Paolucci et al., 2023). These devices aim to enrich the language with extra resources, which act as gender-neutral alternatives to gendered linguistic elements, and allow for a manifest inclusion of gender identities beyond the masculine-feminine binary (Bradley et al., 2019). Individuals choose to use neologistic devices for themselves as they best fit their gender identity and as an open statement of it, rather than using gender-neutralization strategies, which would instead circumvent it (Gautam, 2021). Such innovative solutions are mostly used within LGBTQ+ communities, over informal channels. However, their use and acceptance are on the rise (Waldendorf, 2023; Rose et al., 2023). While there is no *one-fits-all* approach to gender-inclusive language (Lardelli and Gromann, 2023), neologistic devices have naturally emerged as a response to the demand for a direct solution that deserves attention.

In this work, we focus on the use of neomorphemes in en→it translation, a scenario in which gender-related ambiguities – and, consequently, the need for gender-inclusive solutions – are crucial. Indeed, Italian is characterized by a pervasive gender-marking system, which assigns a gender to each noun and every word syntactically linked

to it, including some verbal forms. Coherently, there have been several proposals of neomorpheme paradigms, which currently coexist and are not yet ultimately codified (Sulis and Gheno, 2022). Such proposals promote the use of specific characters in place of gendered morphemes (e.g., masculine -o and feminine -a, as in “uno scienziato” [M], “una scienziata” [F] – en: a scientist). The proposed neomorphemes range from letters of the Latin alphabet (e.g., ‘u’→unu scienziatu), to typographical symbols (e.g., ‘*’→un* scienziat*), to letters of the international phonetic alphabet (IPA), like the Schwa neomorpheme paradigm, which uses the IPA letter ‘ə’ for the singular number (unə scienziatə) and ‘ɜ’ for the plural (alcunɜ scienziatɜ – en: a few scientists) (Baiocco et al., 2023).

Gender Inclusivity in NLP So far, research on gender-inclusive neologistic solutions in NLP has been mainly limited to first explorations in monolingual settings, and mostly confined to English neopronouns. In a pioneering effort, Lauscher et al. (2022) discussed the adoption of neopronouns and formulated a list of desiderata to model the use of pronouns in language technologies. They redefined pronouns as an *open class*, i.e., a class which is not fixed and allows for the inclusion of emerging neopronoun paradigms each user may identify with. This is crucial when dealing with such novel and constantly evolving devices.

In the context of generative tasks, several studies highlight the difficulty of LLMs in handling neopronouns in zero-shot settings (Brandl et al., 2022; Hossain et al., 2023; Ovalle et al., 2023a). Ovalle et al. (2023b) identify byte pair encoding tokenization (Sennrich et al., 2016) as a major cause of LLMs’ shortcomings, coherently with Gaido et al. (2021), which observed the same phenomenon in gender bias investigation. Tokenization, paired with the un-fixed nature of innovative gender-inclusive solutions, may represent a crucial problem for LLMs in correctly generating neomorphemes as well. Indeed, as mentioned above, the range of characters used as neomorphemes is wide and *i*) not all characters are necessarily represented in the training data of LLMs; *ii*) the use of different characters in place of more common gendered morphemes may result in different tokenizations for otherwise identical terms, which in turn could interfere with LLMs’ ability to generate fluent text.

In MT research, the sole experiment in developing systems partially compatible with neologistic

devices is a proof-of-concept built by Saunders et al. (2020) in a gender bias mitigation experiment. They fine-tuned en→de and en→es MT models which use placeholder tags in place of determiners and inflectional morphemes, to be replaced with non-binary forms post-inference. In a broader analysis of gender bias in LLMs, Vanmassenhove (2024) reports that ChatGPT never produces gender-inclusive neomorphemes when translating ambiguous English sentences into Italian, although without specifically prompting the model to do so. The sole analysis dedicated to the use of neologistic devices in MT is the one by Lauscher et al. (2023), which shows how commercial systems fail to deal with English neopronouns, resulting in either misgendering or low-quality outputs.

A major bottleneck hindering the exploration of gender-inclusive neologistic devices in MT is the lack of publicly available evaluation resources. To bridge this gap, in the next section we introduce a dedicated resource: NEO-GATE.

3 The NEO-GATE benchmark

NEO-GATE is designed to evaluate the use of neomorpheme paradigms in en→it translation. Following Lauscher et al. (2022), and extending their desiderata to gender-inclusive translation, we treat neomorphemes as an open class embracing all possible neomorpheme paradigms. To this aim, we design NEO-GATE to be adjustable to any neomorpheme paradigm in Italian, thanks to a set of adaptable references and annotations.

NEO-GATE is built upon GATE (Rarrick et al., 2023), a benchmark for the evaluation of gender bias in MT. In GATE, the gender of human entities is unknown, i.e. there are no elements providing gender information about human referents in the (English) source sentences. GATE also provides target language references which only differ in the feminine/masculine gendered words that refer to human entities (see Table 2). Since in our gender-inclusive translation task we envision the use of neomorphemes for human referents whose gender is unknown, GATE is an ideal candidate corpus as a basis for the creation of our resource. NEO-GATE includes GATE’s test set entries,³ with the addition of references and (word-level) annotations based on a set of placeholder tags, which can be automatically replaced with the de-

³Except for two of GATE’s entries, which do not feature gender-marked terms in the references.

GATE	Source	The department chair said they might hire new professors
	Ref. Masc.	Il direttore del dipartimento ha detto che potrebbero assumere nuovi professori
	Ref. Fem.	La direttrice del dipartimento ha detto che potrebbero assumere nuove professoresse
NEO-GATE	Ref. tagged	<DARTS> direttore<ENDS> del dipartimento ha detto che potrebbero assumere nuov<ENDP> professor<ENDP>
	Annotation	il la <DARTS>; direttore direttrice direttor<ENDS>; nuovi nuove nuov<ENDP>; professori professoresse professor<ENDP>;
NEO-GATE ADAPTED *	Reference	L* direttore* del dipartimento ha detto che potrebbero assumere nuov* professor*
	Annotation	il la l*; direttore direttrice direttor*; nuovi nuove nuov*; professori professoresse professor*;
NEO-GATE ADAPTED ə/3	Reference	Lə direttoreə del dipartimento ha detto che potrebbero assumere nuov3 professor3
	Annotation	il la lə; direttore direttrice direttorə; nuovi nuove nuov3; professori professoresse professor3;

Table 2: Examples of a single entry in GATE, NEO-GATE, and adapted to the two neomorpheme paradigms used in our experiments (§4.2). The terms of interest for our evaluation are highlighted.

sired forms. The tagged references and the annotations are discussed in §3.1, while NEO-GATE’s evaluation metrics are described in §3.2.

3.1 Tagged references and annotations

For each entry in GATE’s test set (see Table 2 for an example), we want to create an additional reference translation featuring neomorphemes. To this aim, for each gendered target word we replace gendered morphemes and function words (articles, prepositions, etc.) with placeholder tags. The placeholders serve to identify words of interest for our task and make this reference adjustable to any neomorpheme paradigm by automatically replacing them with the desired forms. The tagset was designed to cover all parts of the grammar which express grammatical gender, and accounts for distinct singular and plural forms (e.g., the tags <DARTS> and <DARTP> for the singular and plural definite articles respectively). This enables the evaluation of neomorpheme paradigms that use different characters for the singular and the plural case, e.g., the ‘Schwa’ paradigm mentioned in §2. While for content words we only replace the inflectional morpheme with a tag (either <ENDS> or <ENDP>), for function words we use placeholders that cover the whole word. We do so because: *i*) in Italian, some function words are not morphologically derived but paradigmatically opposed (e.g., the definite article singular masculine forms ‘il’ and ‘lo’ vs the feminine form ‘la’); *ii*) as neomorpheme use is not yet settled, there are instances where competing forms exist for a single word and differ in the root part (e.g., the forms ‘l3’ and ‘ə’ have been proposed⁴ for the plural definite article).

⁴The first form was proposed in <https://italianoinclusivo.it/scrittura/>, and the second in <https://effequ.it/schwa/>

Since it would be impossible to account for all existing forms with the sole inflectional placeholders, we replace all function words entirely with dedicated tags (see Appendix B for further details).

We performed the same annotation on a subset of GATE’s dev set as well, so as to have a pool of exemplar sentences for our experiments (see §4). Table 8 in Appendix A describes all the tags used in NEO-GATE, as well as the forms we used to replace them in our experiments. NEO-GATE statistics are reported in Table 3.

	Entries	Tags	Content	Function	Singular	Plural
Test	841	2,479	1,539	940	1,316	1,163
Dev	100	345	211	134	184	161

Table 3: Statistics of NEO-GATE’s test and dev sets.

To ensure the quality of our resource, the references were manually annotated by a linguist following dedicated guidelines.⁵ Using the same guidelines, a second linguist⁶ independently re-annotated a 15% randomly selected subset of target language sentences. Inter-annotator agreement computed with Cohen’s kappa (Cohen, 1960)⁷ on label assignment for the placeholder tags amounts to 0.94, indicating almost perfect agreement (Landis and Koch, 1977). The few disagreements were overlooks and were thus reconciled.

NEO-GATE’s set of annotated words is automatically extracted by comparing the masculine, feminine, and tagged references. It serves to define the words upon which the evaluation is based. It includes the three forms required for the evaluation, i.e the masculine and feminine forms, and the forms with the placeholder tags, which are to be

⁵The guidelines are available in NEO-GATE’s release page.

⁶Both linguists are authors of the paper.

⁷We use scikit-learn (Pedregosa et al., 2011).

replaced with a neomorpheme (e.g., *direttore*, *direttrice*, and *direttore** in ‘NEO-GATE ADAPTED *’, in Table 2).

3.2 Evaluation metrics

While holistic metrics like BLEU (Papineni et al., 2002) have been previously explored to inform the evaluation of gender bias in MT (Bentivogli et al., 2020; Currey et al., 2022), these metrics are not designed to provide fine-grained assessments for specific linguistic phenomena. Rather, they offer a coarse-grained indication of overall translation quality, thus motivating the use of dedicated metrics that allow for pinpointed evaluations, isolating gender from other factors that could impact generic performance. To this aim, we rely on NEO-GATE’s annotations associated with each source sentence (e.g. “*the department chair said they might hire new professors*” in Table 2). Every annotation comprises three forms for each gender-related word: masculine, feminine, and the form with neomorphemes (e.g. “*il la l**”, “*direttore direttrice direttore**”, “*nuovi nuove nuov**”, “*professori professoresse professor**”). In the description of our metrics, we refer to the total number of annotated triplets as ‘*annotations*’ (4 triplets in our example). Scores computation is carried out by scanning the models’ output translations word by word, and checking whether such words match any of the three forms in the annotated triplets. Each matched word increases the ‘*matched*’ count. If the matched form is the one with neomorphemes (e.g. “*direttore**”), we count it as ‘*correct*’. To further monitor models’ behavior we also count the generated words that include a neomorpheme, regardless of their presence in the annotations, in the *found neomorphemes* tally. With these parameters, we compute the following metrics.

Coverage (COV) and accuracy (ACC). As our primary evaluation method, we draw from the metrics defined by Gaido et al. (2020) in the context of binary gender translation. Such a method first computes *coverage* as the ratio of annotated words *matched* in the outputs over NEO-GATE’s *annotations*: $COV = \frac{matched}{annotations}$. This score serves two purposes: *i*) it is indicative of the informativeness of the accuracy evaluation, as a low coverage indicates that the accuracy score described below is calculated over a relatively low number of annotations; *ii*) it can function as an indirect indicator of translation quality (Savoldi et

al., 2022), i.e. a higher coverage suggests that the model generates the expected target words.

On this basis, we then compute *accuracy* as the proportion of *correct* neomorphemes generated by the model over the total number of annotations *matched* in the outputs: $ACC = \frac{correct}{matched}$. This score measures models’ ability to correctly produce neomorphemes.

The combination of these two metrics allows to distinguish between the generation of an annotated word (regardless of its gender) from its gender realization (fem./mas./neom.), thus ensuring pinpointed analyses.

Coverage-weighted accuracy (CWA). For a comprehensive view of models’ overall performance, CWA takes into account both how accurately a model generates neomorphemes and the proportion of *annotations* covered by the evaluation: $CWA = \frac{correct}{matched} * \frac{matched}{annotations}$. This score allows for the comparison of different systems, for which both coverage and accuracy should be taken into account. Indeed, a system’s high accuracy may be the result of an evaluation based on a particularly small set of matched annotations, impairing the comparison with other systems’ performance evaluated on bigger portions of the corpus. While the other metrics serve to investigate each model’s behavior, coverage-weighted accuracy allows for a fairer comparison of different systems.

Mis-generation (MIS). We also consider the case where models generate neomorphemes inappropriately, for instance by applying the use of neomorphemes to words that do not refer to human entities (e.g., *table*→*tavol** instead of ‘*tavolo*’). Such scenario is crucial, as overgeneralizing the use of neomorphemes compromises the intelligibility of the translation. Thus, to we quantify *mis-generations*, i.e. the number of output words – which are not annotated in NEO-GATE – that feature neomorphemes. Accordingly, $MIS = \frac{found\ neomorphemes - correct}{annotations}$. This score complements the evaluation, as it can signal undesired behaviors even despite good accuracy and coverage.

4 Experimental settings

4.1 Models

We experiment with three open, decoder only LLMs. TowerInstruct-7B-v0.1,⁸ is fine-tuned

⁸<https://huggingface.co/Unbabel/TowerInstruct-7B-v0.1>

	BLEU	chrF	TER ↓	BERTSc.	COMET
OPUS-MT	27.53	57.61	58.95	87.42	82.68
GPT-4	32.34	61.11	54.87	88.76	87.05
Tower	<u>30.88</u>	<u>59.41</u>	<u>56.96</u>	<u>88.17</u>	<u>86.21</u>
Mixtral	<u>29.63</u>	<u>58.68</u>	59.35	<u>87.81</u>	<u>86.11</u>
LLama 2	26.28	55.92	61.98	87.02	<u>84.23</u>

Table 4: Translation quality results. Best scores are in **bold**. Cases where LLMs outperform the MT model are underlined.

for MT, whereas the other two – Mixtral-8x7B-Instruct-v0.1 (Jiang et al., 2024) and LLama 2 70B chat (Touvron et al., 2023) – are not specialized for MT. We also include the commercial model GPT-4 (Achiam et al., 2023),⁹ which proved to perform well in gender-inclusive MT experiments (Savoldi et al., 2024). We use the models’ default settings, except for the temperature parameter, which we set to 0 following Peng et al. (2023). We do not include neural MT models as no model currently supports neomorphemes and no dedicated training or fine-tuning data is available.

To ensure the suitability of the selected LLMs for translation-related tasks, we preemptively test their generic en→it translation performance on the FLORES 101 benchmark (Goyal et al., 2022). We prompt the models to translate with a few-shot prompt (see Appendix C). In these experiments, we include opus-mt-en-it,¹⁰ a state-of-the-art neural MT model, as a reference system for translation quality. For general MT evaluation, we use BLEU, chrF (Popović, 2015), TER (Snover et al., 2006), BERTScore (Zhang et al., 2019), and COMET (Rei et al., 2020). Using these metrics allows for a comparative evaluation of translation performance based on different aspects, namely the surface similarity to human-made reference translations (BLEU, chrF, TER), and the semantic adherence to those references (BERTScore) and to the source (COMET). The results in Table 4 show that the LLMs perform very well in MT, often outperforming the SOTA MT system in this setting.

4.2 Neomorphemes

We focus on the two most popular Italian neomorpheme paradigms (Comandini, 2021): *i*) the *Assterisk*, which uses the symbol ‘*’ as a graphemic device in place of regular inflectional morphemes (Haralambous and Dichy, 2018); *ii*) the *Schwa*, which features both a singular form, for which the

character ‘ə’ is used, and a plural form, represented with the character ‘3’ (Baiocco et al., 2023).

For each paradigm, we create a tagset mapping (see Appendix A) that associates the tags used in the tagged references with the desired form for that specific paradigm. As no complete codification of the use and the orthography of neomorphemes in Italian is available (Thornton, 2020), we referenced established resources such as the website *Italiano Inclusivo*,¹¹ and examples found in scientific literature, such as Rosola et al. (2023). As these sources do not cover the whole set of possibly gendered elements in the grammar, we derived the missing forms by analogy from elements of the same class. For example, since none of these sources describes the full set of articulated prepositions, which express gender in Italian, we used the given examples as a model for the rest of the class.

4.3 Prompts

We experiment with one zero-shot and three few-shot formats, illustrated by the examples in Table 5. The few-shot prompts follow the format used in Sánchez et al. (2023), which was found to be useful for controlling gender expression in translation. We instantiate different conceptualizations of the task, ranging from a simple pairing of source sentences directly with gender-inclusive translations, to a ternary opposition of masculine, feminine, and gender-inclusive translations:

◇**Zero-Shot:** a verbalized description of the task is provided without any demonstration.

◇**Direct:** the same verbalized instruction is provided along with demonstrations that include the English source sentence and the gender-inclusive Italian translation.

◇**Binary:** an intermediate gendered (masculine) Italian translation is also included in the format. This format follows the one used in Savoldi et al. (2024), which frames the task as a double output translation. The models are asked to produce a gendered translation first and then a second one with neomorphemes, which should be identical to the first except for the words expressing gender.

◇**Ternary:** two intermediate gendered translations (one masculine, one feminine) are included. The rationale for this format is that, by instantiating a ternary opposition, the models may better identify parts of the target language sentences that should be identical among the three transla-

⁹Model gpt-4-0125-preview

¹⁰<https://huggingface.co/Helsinki-NLP/opus-mt-en-it>

¹¹<https://italianoinclusivo.it/scrittura/>

PROMPT	ROLE	EXAMPLE
Zero-shot	user	Translate the following English sentence into Italian using the neomorpheme ‘*’. To do so, the neomorpheme ‘*’ should be used as a substitute for masculine and feminine morphemes in words that refer to human beings. [English] <{input sentence}> [Italian]
	assistant	<Non compro mai fiori per l* mi* amic*.>
Direct	user	[English] <I never buy flowers for my friends.> [Italian]
	assistant	<Non compro mai fiori per l* mi* amic*.>
Binary	user	[English] <I never buy flowers for my friends.> [Italian, gendered]
	assistant	<Non compro mai fiori per i miei amici.> [Italian, neomorpheme] <Non compro mai fiori per l* mi* amic*.>
Ternary	user	[English] <I never buy flowers for my friends.> [Italian, masculine]
	assistant	<Non compro mai fiori per i miei amici.> [Italian, feminine] <Non compro mai fiori per le mie amiche.> [Italian, neomorpheme] <Non compro mai fiori per l* mi* amic*.>

Table 5: Examples of all the prompts used in our experiments. The few-shots prompt examples include the Asterisk neomorpheme. Words expressing gender are highlighted.

tions and the parts that should differ, i.e. those expressing gender. Framing the task as a triple output translation could help the models infer that the gender expressed in the third translation should be something other than masculine or feminine.

We enclose the exemplar sentences in angle brackets <>. Models are expected to reproduce this structure, thus facilitating the extraction of the final translation from the output in postprocessing.

All four models expect prompts in a ‘chat’ format, with `user` messages providing input and `assistant` messages representing the model’s desired output.¹² For the few-shot prompts we adhere to this structure, whereas for the zero-shot prompts, we only provide a single `user` message.

Demonstrations In the few-shot settings (i.e. Direct, Binary, Ternary) we included 1, 4, and 8 task demonstrations in the prompts. The extremes were chosen as the minimum necessary to elicit in-context learning (1) and a compromise between a high number of demonstrations and the computational cost of inference (8). The exemplar sentences were selected from NEO-GATE’s dev set (§3.1). The exemplars were chosen so as to represent the average tag *density* of the dev set, i.e., the number of tags in each reference, and to offer a balanced mix of singular and plural tags. The prompts were then formatted using each paradigm’s tagset mapping before presenting them to the model.

¹²https://huggingface.co/docs/transformers/main/en/chat_templating

5 Results and discussion

ASTERISK	COV ↑	ACC ↑	CWA ↑	MIS ↓
GPT-4	57.08	74.63	42.60	45.78
Tower	77.57	0.00	0.00	0.00
Mixtral	35.22	37.92	13.35	52.20
LLama 2	56.72	0.57	0.32	16.70
SCHWA	COV ↑	ACC ↑	CWA ↑	MIS ↓
GPT-4	46.91	60.19	28.24	72.77
Tower	77.25	0.00	0.00	0.00
Mixtral	30.05	27.79	8.35	61.44
LLama 2	57.60	0.35	0.20	12.79

Table 6: Zero-shot setting results. We report the coverage (COV), accuracy (ACC), coverage-weighted accuracy (CWA), and mis-generation (MIS) scores.

5.1 Zero-shot results

The results of our zero-shot experiments are reported in Table 6, which unveils very different model behaviors. On the one hand, GPT-4 and Mixtral achieve significantly higher accuracy scores compared to the other two models, with GPT doubling Mixtral’s performance. The accuracy scores indicate that, out of the matched terms, GPT correctly generated 74.63% Asterisk neomorphemes and 60.19% Schwa neomorphemes, with Mixtral reaching 37.92% and 27.79% respectively. Accounting for coverage, the gap widens further, with GPT’s coverage-weighted accuracy amounting to more than three times that of Mixtral (42.60 and 28.24, vs 13.35 and 8.35). Both models tend to

		Asterisk			Schwa		
		Direct	Binary	Ternary	Direct	Binary	Ternary
GPT-4	1-	64.26	71.24	80.11	64.26	64.70	71.00
	4-	69.34	68.46	74.26	72.00	68.54	72.45
	8-	71.68	69.46	74.43	73.17	70.88	71.28
Tower	1-	76.76	76.24	73.62	77.65	73.17	67.16
	4-	78.30	75.23	70.79	76.89	69.30	64.62
	8-	77.85	76.48	73.42	76.28	69.30	63.09
Mixtral	1-	54.22	52.60	52.32	45.06	20.61	23.88
	4-	67.37	61.64	56.03	64.54	50.58	42.23
	8-	72.13	64.86	50.79	70.27	58.17	53.93
LLama 2	1-	54.34	54.26	54.58	62.04	47.08	44.70
	4-	62.44	61.84	59.66	66.84	57.97	52.76
	8-	64.02					

(a) Coverage percentage scores.

		Asterisk			Schwa		
		Direct	Binary	Ternary	Direct	Binary	Ternary
GPT-4	1-	70.62	30.18	4.48	50.53	47.01	24.20
	4-	68.70	80.67	44.27	57.59	81.93	57.24
	8-	67.02	80.49	40.11	58.82	82.93	60.50
Tower	1-	1.89	2.96	6.30	1.35	3.47	13.03
	4-	1.85	5.42	10.48	2.52	10.65	20.79
	8-	2.80	4.85	13.24	3.65	11.06	19.25
Mixtral	1-	56.18	85.66	90.21	39.66	71.23	80.24
	4-	39.04	78.66	87.40	30.25	70.89	78.41
	8-	28.08	73.69	82.02	17.51	62.27	72.03
LLama 2	1-	6.90	6.43	6.43	4.10	8.91	7.40
	4-	1.81	8.74	4.41	1.93	7.03	4.43
	8-	1.95					

(b) Accuracy percentage scores.

Figure 1: Coverage and accuracy results in the few-shot settings. Darker shades indicate better performance.

produce considerable amounts of mis-generations, which are most often higher than the respective coverage scores. This implies that **GPT and Mixtral generate plenty of neomorphemes, but they use them incorrectly for the most part**. Regardless, according to all the metrics, both models perform better with the Asterisk paradigm rather than the Schwa, possibly due to the latter’s use of distinct singular (ə) and plural (ɜ) forms, adding a further challenge to the task.

On the other hand, **LLama 2 and Tower severely under-generate neomorphemes**, regardless of the paradigm. More specifically, LLama’s near zero accuracy scores (0.57 and 0.35) paired with its low mis-generation scores (16.70 and 12.79) indicate that LLama 2 rarely generates neomorphemes and, when it does, it uses them inaccurately. Finally, Tower’s high coverage scores (77.57 and 77.25) combined with the rest of the metrics, all of which report 0 scores, indicate that the model produces fluent, gendered outputs and never generates neomorphemes in the zero-shot setting. This can be due to the fact that in Tower-instruct’s fine-tuning data set, TowerBlocks,¹³ our neomorpheme characters are practically absent (3 occurrences of ‘ɜ’ in English segments, and no occurrences at all of ‘ɜ’ and ‘*’). However, since the development data of the other two models is not publicly available, we cannot further investigate

¹³<https://huggingface.co/datasets/Unbabel/TowerBlocks-v0.1/>

this hypothesis and draw definitive conclusions.

5.2 Few-shots experiments

For the few-shot experiments we report each of the four metrics separately. We do not report all LLama 2 scores because in some cases, namely all the 8-shots settings, the model struggled to reproduce the format described in §4. In such instances, LLama 2 failed to insert the angle brackets or the labels we included in our prompts, and its outputs contained too many hallucinations to be automatically post-processed and evaluated. As the model did not seem to yield better performance or exhibit interesting phenomena in 8-shot settings, the additional effort required to process its unpredictable outputs was unjustified. Therefore, we only report the scores of one of the 8-shots settings outputs, namely the Asterisk, Direct prompt setting.

5.2.1 Coverage and accuracy

The coverage and accuracy scores are reported in Figure 1. **Looking at coverage (1a), we observe that few-shot prompting generally leads to improvements compared to the zero-shot results. Also, for Mixtral and LLama, the scores increase at higher numbers of demonstrations.** As for the prompts, the Direct format generally produces higher coverage scores, with only GPT performing better with the Ternary format. Interestingly, the neomorpheme paradigm has an impact on coverage, as we see generally higher scores with the Asterisk paradigm compared to the

		Asterisk			Schwa		
		Direct	Binary	Ternary	Direct	Binary	Ternary
GPT-4	1 -	45.38	21.50	3.59	32.47	30.42	17.18
	4 -	47.64	55.22	32.88	41.47	56.15	41.47
	8 -	48.04	55.91	29.85	43.04	58.78	43.11
Tower	1 -	1.45	2.26	4.64	1.05	2.54	8.75
	4 -	1.45	4.08	7.42	1.94	7.38	13.44
	8 -	2.18	3.71	9.72	2.78	7.66	12.14
Mixtral	1 -	30.46	45.06	47.20	17.87	14.68	19.16
	4 -	26.30	48.48	48.97	19.52	35.86	33.12
	8 -	20.25	47.80	41.66	12.30	36.22	38.85
LLama 2	1 -	3.75	3.49	3.51	2.54	4.19	3.31
	4 -	1.13	5.40	2.63	1.29	4.08	2.34
	8 -	1.25					

Figure 2: Coverage-weighted accuracy percentage scores for the few-shot settings. Darker shades indicate better performance.

Schwa. As discussed in §5.2.3, this can be ascribed to the models’ tendency to produce more mis-generations with the latter.

Coverage, however, is only informative of the proportion of annotated terms the models generated and disregards how many of those words include neomorphemes. **Looking at accuracy (1b) we find that all models improve their performance in at least one setting** compared to the zero-shot experiments, confirming the benefits of in-context learning for generative tasks involving neologistic expressions (Hossain et al., 2023). Mixtral and GPT are confirmed as the models which produce the highest rates of correct neomorphemes, with the first topping at 90.21 and the latter at 82.93 accuracy. On the contrary, Tower and LLama 2 are unfit for this task despite their improvements, as their scores remain low.

Surprisingly, **a greater number of demonstrations does not necessarily lead to higher accuracy.** While coverage generally increased with more demonstrations, this trend generally holds true for accuracy only for GPT and Tower, indicating that they generate more neomorphemes and do so more accurately. On the contrary, the accuracy of LLama 2 and Mixtral significantly decreases with more demonstrations. Paired with their rising coverage, this indicates that they produce fewer neomorphemes and more gendered terms. Both behaviors may result from systems better modeling the task with more demonstrations, as LLama’s

		Asterisk			Schwa		
		Direct	Binary	Ternary	Direct	Binary	Ternary
GPT-4	1 -	46.51	27.59	5.24	43.24	50.26	28.04
	4 -	33.52	53.05	20.94	29.77	57.81	31.26
	8 -	25.86	44.74	15.45	25.21	46.87	27.11
Tower	1 -	26.46	6.86	11.42	0.81	10.17	26.14
	4 -	6.41	12.46	18.56	2.10	19.81	28.24
	8 -	5.85	7.22	12.51	2.46	18.64	28.12
Mixtral	1 -	52.32	143.81	102.38	61.64	198.55	180.48
	4 -	24.45	90.84	84.31	25.09	104.52	120.29
	8 -	13.72	4.60	58.53	9.96	57.56	60.51
LLama 2	1 -	36.30	34.65	35.09	13.03	35.26	40.90
	4 -	10.65	18.48	15.97	4.20	20.73	19.00
	8 -	5.73					

Figure 3: Mis-generation percentage scores for the few-shot settings. Higher scores (darker shades) indicate worse performance.

and Mixtral’s higher accuracy and lower coverage in the 1 shot settings may be due to fortuitous correct generations of neomorphemes in a context where they over-generate them. We discuss this aspect below, looking at mis-generation (§5.2.3).

As for the neomorpheme paradigms, Mixtral and Tower perform better with the Asterisk and the Schwa respectively, as in the zero-shot experiments. GPT does not seem to be consistently affected by the neomorpheme paradigm. LLama presents negligible differences between the paradigms as well. The ability to more correctly generate one neomorpheme over another is possibly due to models’ robustness to likely unseen grammatical paradigms and to the representations of the specific characters used as neomorphemes in each model’s training data. Unfortunately, we cannot investigate this aspect as such data is not publicly available, with the exception of Tower’s fine-tuning data set, as mentioned in §5.1.

5.2.2 Coverage-weighted accuracy

To compare models’ overall performance in gender-inclusive MT, we look at coverage-weighted accuracy in Figure 2. This metric offers a comprehensive view of model performance in each setting, allowing for a fair comparison of different systems in light of both coverage and accuracy.

We first find that **all models improve their performance in the few-shots experiments.** The upside offered by in-context learning is notable, and there is arguably room for improvement at higher

numbers of demonstrations. GPT and Mixtral are confirmed as the best models, and the gap between them narrows significantly with respect to the zero-shot experiments. In the best configurations, GPT scores 58.78 (Schwa, Binary, 8 shots) and Mixtral scores 48.97 (Asterisk, Ternary, 4 shots). Generally, GPT performs better with the Binary prompt and 4 or 8 shots, whereas Mixtral achieves its best results in the Binary/Ternary, 4/8 shots region, especially with the Asterisk paradigm. As for the other two models, despite the very low scores Tower generally outperforms the ten times bigger LLama 2, but both come across as unfit for this task.

5.2.3 Mis-generation

So far, the discussion of our few-shots experiments has focused on the correct generation of neomorphemes when referring to human entities, thus only on relevant phenomena we annotated in our test set. To better investigate models' behavior, we look at mis-generation, i.e. inappropriate neomorphemes generations, as well (Figure 3).

We first note that **Mixtral stands out as the model producing the most mis-generations**, especially in the 1 and 4 shots, Binary and Ternary region. Table 7 reports examples of mis-generation from Mixtral's outputs.

Source	I hope the shaman can help us.
Annotation	lo la lə; sciamano sciamana sciamanə;
Output	Sperə che lə <u>sciamanə</u> possa aiutarci.
Source	They asked everyone to remain silent.
Annotation	tutti tutte tutt3;
Output	Hanno chiesto a t3 di <u>rimanerə</u> in <u>silenziə</u> .

Table 7: Examples of mis-generation found in Mixtral's Schwa, 1 shot, Binary prompt outputs. Words containing neomorphemes are underlined, mis-generations are in bold.

As hypothesized in §5.2.1, in these settings Mixtral over-generates neomorphemes, resulting in both correct generations and mis-generations. This behavior is reflected in the high accuracy and low coverage: by over-generating neomorphemes, Mixtral produces fewer gendered words – which would contribute to coverage – and many words that are either *a)* correct or *b)* mis-generations. With more task demonstrations, Mixtral generates significantly fewer mis-generations, and while its accuracy decreases the coverage improves, meaning that it produces better formed outputs. Mixtral's example testifies to how the mis-generation metric complements the analysis of models' be-

havior, as it sheds light on unwanted phenomena related to neomorphemes usage, which coverage and accuracy alone (or combined) cannot signal.

Similarly to Mixtral, LLama 2 produces more mis-generations given fewer demonstrations and progressively mitigates this behavior when given more. With an opposite trend, GPT generates fewer mis-generations, and Tower even less. However, the best performing settings of both models are also the ones in which they produce the most mis-generations. Hence, future work should focus on improving the ratio of correctly generated neomorphemes over the total neomorphemes generated by these models.

6 Conclusions

We discussed a neologistic approach to gender-inclusive machine translation, an underexplored area constrained by the lack of publicly available dedicated data. Our first contribution, the release of the NEO-GATE benchmark, allowed us to give a first fundamental impulse to research in this direction. As a second contribution, we explored the possibility of performing gender-inclusive translation from English to Italian with four popular Large Language Models: three open models – Mixtral, Tower, and LLama 2 – and a commercial one – GPT-4. Our comparisons across different prompting settings reveal that GPT-4 and Mixtral generally exhibit promising results when properly prompted, while LLama 2 and Tower are unfit for the task. More specifically, models' understanding of the task is significantly influenced by prompt complexity, the number of demonstrations, and the specific characters employed as neomorphemes (possibly depending on the representation of those characters in each model's training data).

While our investigation suggests LLMs' potential for neologistic gender-inclusive MT, there remains room for improving their accuracy. NEO-GATE and the analyses presented herein lay the groundwork for rising to the challenge and for future research on gender-inclusive MT tailored to existing neologistic paradigms, and those that may emerge in this new and evolving landscape.

7 Acknowledgements

We acknowledge the support of the PNRR project FAIR - Future AI Research (PE00000013), under the NRRP MUR program funded by the NextGenerationEU.

References

- Achiam, Josh, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Attig, Remy and Artemis López. 2020. Queer Community Input in Gender-Inclusive Translations. *Linguistic Society of America [Blog]*, June 23.
- Baiocco, Roberto, Fau Rosati, and Jessica Pistella. 2023. Italian proposal for non-binary and inclusive language: The schwa as a non-gender-specific ending. *Journal of Gay & Lesbian Mental Health*, 27(3):248–253, July. Publisher: Routledge eprint: <https://doi.org/10.1080/19359705.2023.2183537>.
- Bentivogli, Luisa, Beatrice Savoldi, et al. 2020. Gender in Danger? Evaluating Speech Translation Technology on the MuST-SHE Corpus. In *58th ACL*, pages 6923–6933.
- Blodgett, Su Lin, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of “bias” in NLP. In Jurafsky, Dan, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, Online, July. Association for Computational Linguistics.
- Bradley, Evan D., Julia Salkind, Ally Moore, and Sofi Teitsort. 2019. Singular ‘they’ and novel pronouns: Gender-neutral, nonbinary, or both? *Proceedings of the Linguistic Society of America*, 4:36:1–7, March.
- Brandl, Stephanie, Ruixiang Cui, and Anders Søgaard. 2022. How conservative are language models? adapting to the introduction of gender-neutral pronouns. In Carpuat, Marine, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3624–3630, Seattle, United States, July. Association for Computational Linguistics.
- Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Cohen, Jacob. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.
- Comandini, Gloria. 2021. Salve a tuttə, tutt*, tuttu, tuttx e tutt@: l’uso delle strategie di neutralizzazione di genere nella comunità queer online. indagine su un corpus di italiano scritto informale sul web. *Testo e Senso*, 23:43–64.
- Currey, Anna, Maria Nadejde, Raghavendra Reddy Pappagari, Mia Mayer, Stanislas Lauly, Xing Niu, Benjamin Hsu, and Georgiana Dinu. 2022. MT-GenEval: A counterfactual and contextual dataset for evaluating gender accuracy in machine translation. In Goldberg, Yoav, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 4287–4299, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.
- Dev, Sunipa, Masoud Monajatipoor, Anaelia Ovalle, Arjun Subramonian, Jeff Phillips, and Kai-Wei Chang. 2021. Harms of gender exclusivity and challenges in non-binary representation in language technologies. In Moens, Marie-Francine, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1968–1994, Online and Punta Cana, Dominican Republic, November. Association for Computational Linguistics.
- Gaido, Marco, Beatrice Savoldi, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2020. Breeding gender-aware direct speech translation systems. In Scott, Donia, Nuria Bel, and Chengqing Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3951–3964, Barcelona, Spain (Online), December. International Committee on Computational Linguistics.
- Gaido, Marco, Beatrice Savoldi, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. How to split: the effect of word segmentation on gender bias in speech translation. In Zong, Chengqing, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3576–3589, Online, August. Association for Computational Linguistics.
- Gautam, Vasundara. 2021. Guest lecture in pronouns: Vagrant. <https://link.medium.com/viFawWyPVHb>. Accessed: Feb 20, 2024.
- Goyal, Naman, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

- Gustafsson Sendén, Marie, Emma A. Bäck, and Anna Lindqvist. 2015. Introducing a gender-neutral pronoun in a natural gender language: the influence of time on attitudes and behavior. *Frontiers in Psychology*, 6.
- Haralambous, Yannis and Joseph Dichy. 2018. Graphemic Methods for Gender-Neutral Writing. In *Graphemics in the 21st Century, Brest 2018*, volume Graphemics in the 21st Century: 2018 Conference, pages 41 – 89, Brest, France, Jun. Fluxus Editions.
- Hossain, Tamanna, Sunipa Dev, and Sameer Singh. 2023. MISGENDERED: Limits of large language models in understanding pronouns. In Rogers, Anna, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5352–5367, Toronto, Canada, July. Association for Computational Linguistics.
- Jiang, Albert Q, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Kaplan, Jennifer Marisa. 2022. Pluri-grammars for pluri-genders: Competing gender systems in the nominal morphology of non-binary french. *Languages*, 7(4).
- Landis, J. Richard and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics*, 33(1):159–174.
- Lardelli, Manuel and Dagmar Gromann. 2023. Gender-fair post-editing: A case study beyond the binary. In Nurminen, Mary, Judith Brenner, Maarit Koponen, Sirkku Latomaa, Mikhail Mikhailov, Frederike Schierl, Tharindu Ranasinghe, Eva Vanmassenhove, Sergi Alvarez Vidal, Nora Aranberri, Mara Nunziatini, Carla Parra Escartún, Mikel Forcada, Maja Popovic, Carolina Scarton, and Helena Moniz, editors, *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 251–260, Tampere, Finland, June. European Association for Machine Translation.
- Lauscher, Anne, Archie Crowley, and Dirk Hovy. 2022. Welcome to the Modern World of Pronouns: Identity-Inclusive Natural Language Processing beyond Gender. In *29th COLING*, pages 1221–1232.
- Lauscher, Anne, Debora Nozza, Ehm Miltersen, Archie Crowley, and Dirk Hovy. 2023. What about “em”? how commercial machine translation fails to handle (neo-)pronouns. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 377–392, Toronto, Canada, July. Association for Computational Linguistics.
- McGaughey, Sebastian. 2020. Understanding neopronouns. *The Gay & Lesbian Review Worldwide*, 27:27+.
- Min, Sewon, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Rethinking the Role of Demonstrations: What Makes In-Context Learning Work? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11048–11064, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.
- Misiek, Szymon. 2020. Misgendered in translation?: Genderqueerness in polish translations of english-language television series. *Anglica. An International Journal of English Studies*, pages 165–185.
- Motschenbacher, Heiko. 2014. Grammatical gender as a challenge for language policy: The (im)possibility of non-heteronormative language use in German versus English. *Language policy*, 13(3):243–261.
- Ovalle, Anaelia, Palash Goyal, Jwala Dhamala, Zachary Jagers, Kai-Wei Chang, Aram Galstyan, Richard Zemel, and Rahul Gupta. 2023a. “i’m fully who i am”: Towards centering transgender and non-binary voices to measure biases in open language generation. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, FAccT ’23, page 1246–1266, New York, NY, USA. Association for Computing Machinery.
- Ovalle, Anaelia, Ninareh Mehrabi, Palash Goyal, Jwala Dhamala, Kai-Wei Chang, Richard Zemel, Aram Galstyan, Yuval Pinter, and Rahul Gupta. 2023b. Are you talking to [ˈxem] or [ˈxˈ,em]? on tokenization and addressing misgendering in llms with pronoun tokenization parity. *arXiv preprint arXiv:2312.11779*.
- Paolucci, Angela Balducci, Manuel Lardelli, and Dagmar Gromann. 2023. Gender-fair language in translation: A case study. In Vanmassenhove, Eva, Beatrice Savoldi, Luisa Bentivogli, Joke Daems, and Janiça Hackenbuchner, editors, *Proceedings of the First Workshop on Gender-Inclusive Translation Technologies*, pages 13–23, Tampere, Finland, June. European Association for Machine Translation.
- Papineni, Kishore, Salim Roukos, et al. 2002. BLEU: a Method for Automatic Evaluation of Machine Translation. In *40th ACL*, pages 311–318.
- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Peng, Keqin, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and

- Dacheng Tao. 2023. Towards making the most of ChatGPT for machine translation. In Bouamor, Houda, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5622–5633, Singapore, December. Association for Computational Linguistics.
- Piazzolla, Silvia Alma, Beatrice Savoldi, and Luisa Bentivogli. 2023. Good, but not always fair: An evaluation of gender bias for three commercial machine translation systems. *HERMES - Journal of Language and Communication in Business*, (63):209–225, Dec.
- Piergentili, Andrea, Dennis Fucci, Beatrice Savoldi, Luisa Bentivogli, and Matteo Negri. 2023a. Gender neutralization for an inclusive machine translation: from theoretical foundations to open challenges. In *Proceedings of the First Workshop on Gender-Inclusive Translation Technologies*, pages 71–83, Tampere, Finland, June. European Association for Machine Translation.
- Piergentili, Andrea, Beatrice Savoldi, Dennis Fucci, Matteo Negri, and Luisa Bentivogli. 2023b. Hi guys or hi folks? benchmarking gender-neutral machine translation with the GeNTE corpus. In Bouamor, Houda, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14124–14140, Singapore, December. Association for Computational Linguistics.
- Popović, Maja. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395.
- Rarrick, Spencer, Ranjita Naik, Varun Mathur, Sundar Poudel, and Vishal Chowdhary. 2023. Gate: A challenge set for gender-ambiguous translation examples. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '23, page 845–854, New York, NY, USA. Association for Computing Machinery.
- Rei, Ricardo, Craig Stewart, Ana C. Farinha, and Alon Lavie. 2020. COMET: A Neural Framework for MT Evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702.
- Rose, Ell, Max Winig, Jasper Nash, Kyra Roepke, and Kirby Conrod. 2023. Variation in acceptability of neologistic English pronouns. *Proceedings of the Linguistic Society of America*, 8(1):5526, April.
- Rosola, Martina, Simona Frenda, Alessandra Teresa Cignarella, Matteo Pellegrini, Andrea Marra, Mara Floris, et al. 2023. Beyond obscuration and visibility: Thoughts on the different strategies of gender-fair language in Italian. *CLiC-it 2023. Proceedings of the 9th Italian Conference on Computational Linguistics. Venice, Italy, November 30-December 2, 2023.*, 3596:1–10.
- Sánchez, Eduardo, Pierre Andrews, Pontus Stenetorp, Mikel Artetxe, and Marta R Costa-jussà. 2023. Gender-specific machine translation with large language models. *arXiv preprint arXiv:2309.03175*.
- Sarmiento, Miguel Angel. 2015. La e para la desexualización del género en beneficio de la motivación de ele en suecia : Revitalizando la propuesta de Álvaro garcía meseguer. In *La enseñanza de ELE centrada en el alumno .*, pages 863–872. ASELE.
- Saunders, Danielle, Rosie Sallis, and Bill Byrne. 2020. Neural Machine Translation Doesn't Translate Gender Coreference Right Unless You Make It. In *2nd Workshop on Gender Bias in Natural Language Processing*, pages 35–43.
- Savoldi, Beatrice, Marco Gaido, et al. 2021. Gender bias in machine translation. *Transactions of the Association for Computational Linguistics*, 9:845–874.
- Savoldi, Beatrice, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2022. Under the morphosyntactic lens: A multifaceted evaluation of gender bias in speech translation. In Muresan, Smaranda, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1807–1824, Dublin, Ireland, May. Association for Computational Linguistics.
- Savoldi, Beatrice, Andrea Piergentili, Dennis Fucci, Matteo Negri, and Luisa Bentivogli. 2024. A prompt response to the demand for automatic gender-neutral translation. In Graham, Yvette and Matthew Purver, editors, *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 256–267, St. Julian's, Malta, March. Association for Computational Linguistics.
- Sendén, Marie Gustafsson, Emma Renström, and Anna Lindqvist. 2021. Pronouns beyond the binary: The change of attitudes and use over time. *Gender & Society*, 35(4):588–615.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In Erk, Katrin and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August. Association for Computational Linguistics.
- Snover, Matthew, Bonnie Dorr, et al. 2006. A Study of Translation Edit Rate with Targeted Human Annotation. In *7th AMTA*, pages 223–231.
- Sulis, Gigliola and Vera Gheno. 2022. The debate on language and gender in Italy, from the visibility of women to inclusive language (1980s–2020s). *The Italianist*, 42(1):153–183.

Thornton, Anna Maria. 2020. Genere e igiene verbale: l'uso di forme con *oin* italiano | annali del dipartimento di studi letterari, linguistici e comparati. sezione linguistica. *Annali Del Dipartimento Di Studi Letterari, Linguistici E Comparati. Sezione Linguistica*, 11:11–54.

Touvron, Hugo, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Vanmassenhove, Eva. 2024. Gender bias in machine translation and the era of large language models. *arXiv preprint arXiv:2401.10016*.

Waldendorf, Anica. 2023. Words of change: The increase of gender-inclusive language in German media. *European Sociological Review*, September.

Zhang, Tianyi, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2019. BERTScore: Evaluating Text Generation with BERT. In *International Conference on Learning Representations*.

A Tagset and annotation

Table 8 reports the complete tagset used in NEO-GATE, as well as the tagset mappings for the Schwa and the Asterisk paradigms.

B Function words anchoring

We include an additional information for function words which maps them to an anchor, in respect to which they are expected to be correctly positioned. This check allows for a more precise evaluation of function words, as it ensures that the evaluation is performed on the appropriate function word, and not on other ones which may occur in the sentence.

An anchor consists in the longest possible sub-word common to the masculine, feminine, and tagged content word which the function word is syntactically linked to. Looking at Table 9, the first Annotation reports an example of anchor for a function word: ‘student=1’. It indicates that the sub-word sequence ‘student’ is the anchor for the function word forms ‘il la l*’, meaning that if one of the three forms is found it will only be evaluated if the anchor is found immediately after it (i.e., at a distance of 1 word). Similarly, the second annotation of the table reports anchor annotations for two function words. The first, ‘amic=2’ indicates that if one of the three forms ‘i le l*’ is found, it will only be evaluated if the anchor ‘amic’ is found at a distance of two words. The second anchor annotation ‘amic=1’ maps the function word forms ‘tuo

tue tu*’ to the same anchor ‘amic’, which should be positioned one word after.

We did not include anchor annotations in the main body to simplify the examples. However, all function words annotated in NEO-GATE are assigned with anchors, including the ones reported in the examples throughout the paper.

C Translation experiments prompt

Table 10 reports the prompt we used to assess the general translation quality of the systems, as discussed in §4.1. We include three demonstrations taken from FLORES’ dev set, so as to provide the LLMs with an interaction structure to reproduce. This facilitates the process of filtering out extra comments and hallucination produced by the models, and extract the output translation.

TAG	Description	Masculine	Feminine	Asterisk	Schwa
<ENDS>	inflectional morpheme (word ending), singular	o, e, tore	a, essa, trice	*	ə
<ENDP>	inflectional morpheme (word ending), plural	i, tori	e, esse, trici	*	ɜ
<DARTS>	definite article, singular	il, lo, l'	la, l'	l*	lə
<DARTP>	definite article, plural	i, gli	le	l*	lɜ
<IART>	indefinite article	uno, un	una, un'	un*	unə
<PARTP>	partitive article, plural	dei, degli	delle	de*	deɜ
<PREPdiS>	articulated preposition with root 'di', singular	del, dello, dell'	della, dell'	dell*	dellə
<PREPdiP>	articulated preposition with root 'di', plural	dei, degli	delle	dell*	dellɜ
<PREPaS>	articulated preposition with root 'a', singular	al, allo, all'	alla, all'	all*	allə
<PREPaP>	articulated preposition with root 'a', plural	agli, ai	alle	all*	allɜ
<PREPdaS>	articulated preposition with root 'da', singular	dal, dallo, dall'	dalla, dall'	dall*	dallə
<PREPdaP>	articulated preposition with root 'da', plural	dagli	dalle	dall*	dallɜ
<PREPinP>	articulated preposition with root 'in', plural	negli	nelle	nell*	nellɜ
<PREPsuS>	articulated preposition with root 'su', singular	sul, sullo, sull'	sulla, sull'	sull*	sullə
<PREPsuP>	articulated preposition with root 'su', plural	sugli	sulle	sull*	sullɜ
<DADJquelS>	demonstrative adjective (far), singular	quel, quello, quell'	quella, quell'	quell*	quellə
<DADJquelP>	demonstrative adjective (far), plural	quegli	quelle	quell*	quellɜ
<DADJquestS>	demonstrative adjective (near), singular	questo, quest'	questa, quest'	quest*	questə
<DADJquestP>	demonstrative adjective (near), plural	questi	queste	quest*	questɜ
<POSS1S>	possessive adjective, 1st person singular, singular	mio	mia	mi*	miə
<POSS1P>	possessive adjective, 1st person singular, plural	miei	mie	mi*	miɜ
<POSS2S>	possessive adjective, 2nd person singular, singular	tuo	tua	tu*	tuə
<POSS2P>	possessive adjective, 2nd person singular, plural	tuo	tue	tu*	tuɜ
<POSS3S>	possessive adjective, 3rd person singular, singular	suo	sua	su*	suə
<POSS3P>	possessive adjective, 3rd person singular, plural	suoi	sue	su*	suɜ
<POSS4S>	possessive adjective, 1st person plural, singular	nostro	nostra	nostr*	nostrə
<POSS4P>	possessive adjective, 1st person plural, plural	nostri	nostre	nostr*	nostrɜ
<PRONDOBJS>	direct object pronoun, singular	lo	la	l*	lə
<PRONDOBJP>	direct object pronoun, plural	li	le	l*	lɜ

Table 8: The full tagset used in NEO-GATE and the tagset mappings to the Italian gendered forms and the desired forms in the Asterisk and Schwa nomorpheme paradigms.

Source	The student was worried about going off topic.
Tagged reference	L* student* era preoccupat* di andare fuori tema.
Annotation	lo la l* student=1 ; studente studentessa student*; preoccupato preoccupata preoccupat*;
Source	Come out to the balcony and let your friends see you.
Tagged reference	Esci sul balcone e lascia che l* tu* amic* ti vedano.
Annotation	i le l* amic=2 ; tuoi tue tu* amic=1 ; amici amiche amic*;

Table 9: Examples of NEO-GATE annotations including anchors. The Tagged references and the Annotations are adapted to the Asterisk paradigm, the anchors are in bold.

Role	Content
user	Translate the following English sentence into Italian: <Construction is ongoing for five new skyscrapers at the site, with a transportation center and memorial park in the middle.>
assistant	<Nel sito sono in corso i lavori di costruzione di cinque nuovi grattacieli, con un centro trasporti e un parco alla memoria nel mezzo.>
user	<The game is based on the Second Battle of Fallujah, a vicious battle between American and Iraqi forces.>
assistant	<Il videogioco simula la seconda battaglia di Fallujah, un conflitto cruento tra le forze armate americane e quelle irachene.>
user	<Born in the Croatian capital, Zagreb, Bobek gained fame while playing for Partizan Belgrade.>
assistant	<Nato nella capitale croata Zagabria, Bobek è diventato famoso giocando nel Partizan Belgrado.>
user	<{input sentence}>

Table 10: The 3 shots prompt used in the general translation preliminary experiments.

Research: Translators & Users

Prompting ChatGPT for Translation: A Comparative Analysis of Translation Brief and Persona Prompts

Sui He

School of Culture and Communication
Swansea University
United Kingdom
sui.he@swansea.ac.uk

Abstract

Prompt engineering has shown potential for improving translation quality in LLMs. However, the possibility of using translation concepts in prompt design remains largely underexplored. Against this backdrop, the current paper discusses the effectiveness of incorporating the conceptual tool of “translation brief” and the personas of “translator” and “author” into prompt design for translation tasks in ChatGPT. Findings suggest that, although certain elements are constructive in facilitating human-to-human communication for translation tasks, their effectiveness is limited for improving translation quality in ChatGPT. This accentuates the need for explorative research on how translation theorists and practitioners can develop the current set of conceptual tools rooted in the human-to-human communication paradigm for translation purposes in this emerging workflow involving human-machine interaction, and how translation concepts developed in translation studies can inform the training of GPT models for translation tasks.

1 Introduction

Translation quality is a pivotal topic in the field of machine translation. The development of Large Language Models (LLMs) and the popularization of ChatGPT since its public launch in November 2022 have attracted scholarly interests in im-

proving the quality of translation outputs generated by LLMs. Efforts to improve the quality of these translations have involved both fine-tuning and prompt engineering. Despite these efforts, the performance of popular LLMs in executing translation tasks remains suboptimal, particularly when compared with professional translations used in the language service industry (Jiao et al., 2023). Therefore, the task of enhancing the performance of LLMs in conducting translation tasks continues to be an ongoing effort.

Compared to fine-tuning, prompt engineering provides greater accessibility for ordinary users with translation needs, especially those who operate on the user interfaces of LLMs such as ChatGPT. Most research on prompt engineering for translation purposes draws on concepts such as zero-shot learning rooted in Natural Language Processing (NLP) by feeding sample translations in the context window. In comparison, the possibilities for integrating translation concepts and strategies have received little attention.

From the perspective of advancing translation studies, consolidating the synergy between humans and machines in achieving translation goals at a professional level is crucial. As Lee (2023) rightly notes, “translation as an event can no longer be restricted to translating as an act, given that AI and other communicative modalities will increasingly be drawn into and embedded within the workflow.” For the development of translation research, since most translation concepts are anchored in human-to-human communication, it becomes essential to evaluate their efficacy in the emerging workflow with human-machine communication involved, thereby strengthening the disciplinary foundation of translation studies in this novel context. For translation practice, enhancing

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

our understanding of prompt engineering for translation could inform the approach we take for translator training in the changing landscape. As highlighted in a recent work trend report by Microsoft (2023), 82% of leaders from various sectors stated that their employees will need new competencies – such as AI delegation via prompts – to prepare for the expansion of AI.

In the background, this research investigates the effectiveness of incorporating the notion of “translation brief” and the translator/author dichotomy into prompt design, as an attempt to explore the potential of using conceptual tools rooted in translation studies for improving the quality of LLM-generated translations. In this study, ChatGPT is chosen for its popularity among general users and its user-friendly interface that accommodates individuals with limited computing expertise. Based on two sets of experiments, this research seeks to answer two questions specific to the scope of the current study: 1) Compared to a basic translation command, does a prompt containing information included in a typical translation brief help improve the quality of translation outputs? 2) Drawing on the persona feature of ChatGPT, does assigning the role of “translator” make a difference to the translation quality, with the basic instruction and the role of “author” as reference points?

2 Literature Review

In the guidance for prompt design in ChatGPT published by OpenAI,¹ six strategies are listed for creating effective prompts. Even though there is an improvement in the content of this guideline when compared to its earlier version where only three generic strategies (i.e., show and tell, provide quality data, check your settings; accessed April 2023) were suggested, OpenAI has not yet published any specialized guidance on prompt design for translation purposes in ChatGPT. Nonetheless, scholarly efforts have been made to address this issue.

In the literature, most of the research focuses on prompting GPT models or other LLMs through APIs (Vilar et al., 2023; Hendy et al., 2023; Zhang et al., 2023; Zhu et al., 2023). However, a small number of studies have also explored prompt engineering for translation tasks specifically through the user interface of ChatGPT, drawing on different linguistic concepts. Within this niche area,

¹<https://platform.openai.com/docs/guides/completion/prompt-design>

two main threads have emerged: one is centered around specific translation problems, and the other features a more holistic approach.

Starting with those targeting specific translation problems, Gu (2023) is the only one in the literature to date. Drawing on the default model (GPT-3.5) of ChatGPT, the author utilizes the “in context learning” capability of ChatGPT (i.e., remembering what has been mentioned in the chat) to “teach” it how to translate attributive clauses. Specifically, a translation strategy commonly adopted by translators to render attributive clauses from Japanese into Chinese was used by the author to design a set of prompts: “What is the noun modified by the attributive clause in the following sentence?”, “Place the noun modified by the attributive clause in the subject position of the attributive clause. And then separate [SOURCE SENTENCE] into two sentences”, and finally “Translate the following sentence to Chinese: [SEPARATED SOURCE SENTENCE]”. Although this prescriptive application of a standalone translation strategy fails to take into consideration the dynamic context of handling attributive clauses, this paper presents a very interesting attempt to bring translation strategies into the horizon of prompt engineering.

Turning to the literature which investigates translation at a contextual level, key concepts tested in this group include “domain”, “task”, “part of speech”, “discourse”, and “pivot language” – all of them are well-established topics in translation studies but they have been used in a rather ambiguous way in these works. For instance, Peng et al. (2023) propose the concept of “task-specific prompts” (i.e., “you are a machine translation system”) in their experiment, without concrete instructions on what to expect from a so-called “machine translation system”. The rationale behind this design, according to the authors, rests in the assumption that ChatGPT has been fine-tuned as a conversation system instead of a machine translation system, and this might have limited the translation ability of ChatGPT. Nonetheless, the effectiveness of altering a fine-tuned chatbot into a machine translation system with a single prompt line in the user interface remains questionable. Additionally, the authors test the efficacy of “domain-specific prompts” (e.g., information about the topic or genre of the ST, such as bio-medical or news-style) by providing ChatGPT with both right and

wrong domain information of the ST. This design of using wrong domain information, from the perspective of translation studies, requires careful justification. The results, measured via automated machine translation quality evaluation metrics, suggest that providing task and correct domain information can indeed enhance ChatGPT's translation performance.

Another case in point is Gao et al. (2023). The authors introduce language direction, domain information, and part-of-speech information to their prompt design. Similar to the definition of “domain” in Peng et al. (2023), the authors include information about genre (i.e., news, e-commerce, social, and conversational texts) in their tests. These prompts were run through five different settings to test their efficacy. The results from automatic metrics further validate the usefulness of domain-related information in prompt engineering for translation tasks. Notably, although the outcome of introducing part-of-speech information in prompts was not promising, it suggests an intention to include grammatical segmentations into prompt design, which echoes the problem-oriented approach to enhancing translation quality, as mentioned above in Gu (2023). An interesting observation made by the authors regarding language direction lies in the disparity between high-resource languages and low-resource languages: domain information appears to enhance machine translation quality for high-resource languages but fails to demonstrate a comparable impact on low-resource languages.

To understand the issue related to high versus low resource languages, Jiao et al. (2023) propose a strategy called “pivot prompting”. This notion, bearing similarities to the concept of relay translation, involves instructing ChatGPT to translate the ST into a high-resource language prior to translating it into the target language. Even though the basic prompts were generated by ChatGPT itself without further tweaks, the idea of relay translation turned out to be useful in improving translation quality between distant languages, as the results reported by the authors suggest.

Regarding the topic of context and discourse in translation, whilst all studies mentioned above focus on prompt design for translation at the level of single sentences or small sentence clusters, Wang et al. (2023) take a step forward to the document level. They put forward the concept of

“discourse-aware prompts”, introducing discourse as an evaluation criterion for assessing the quality of prompts in ChatGPT. To identify the best discourse-aware prompt, the authors evaluate a set of basic prompts generated by ChatGPT with two discourse-oriented metrics: one focuses on terminology consistency and another on the accuracy of zero pronoun translation. As can be seen from the design, discourse here is used in its micro sense as document-level coherence. Macro discursal information, such as the function of the ST and target audience, is not taken into consideration when designing the prompts.

The most relevant research to date, drawing on a contextualized approach inspired by translation concepts, is reported by Yamada (2024). There are two sets of experiments in this research. First, the author adopts two concepts – purpose of the translation and target readers – for prompting ChatGPT (GPT-4) to translate, simulating a real-life translation commission for ChatGPT. Instead of providing information about the purpose and target readers, the author designed a prompt that asks ChatGPT to find the information itself: “Translate the following Japanese [source text] into English. Please fulfil the following conditions when translating. Purpose of the translation: *You need to fill in.* Target audience: *You need to fill in.* [source text] *You need to fill in.*” In the segments shown in italics, the author specifies the information that ChatGPT needs to fill in before generating the translation. Second, the concept of dynamic equivalence is utilized, feeding into ChatGPT as a translation strategy alongside a sample translation of a different source text through in-context learning. This combined approach complicates the task of determining whether the concept of “dynamic equivalence” or its illustrative examples play a more significant role in the efficacy of the prompt. To assess the overall effectiveness of this prompt, the author uses cosine similarity of vectors as indicators for semantic proximity and a detailed qualitative evaluation conducted by the author himself, with reference translations generated by DeepL, Google Translate, and ChatGPT (with default prompt “Translate to English”). The author reports that “incorporating the purpose and target readers into prompts indeed altered the generated translations” and that “this transformation [...] generally improved the translation quality by industry standards”. This research features a very

interesting attempt to “teach” ChatGPT to “think” and “act” like a translator via prompts, revealing the potential for training ChatGPT with knowledge generated by translation scholars.

Overall, the current landscape of prompt design in ChatGPT features important attempts to enhance its capability in executing translation tasks. However, a critical issue with these endeavors lies in the fact that the concepts being used in the prompts (e.g., “news-style”) are too general to be informative, and some of the approaches (e.g., the out-of-context application of prescriptive translation strategies) bear striking resemblances to what happened in the early days in translation studies. The design of prompts shows that these research efforts have touched upon some key conceptual tools for translation, revealing the potential benefit that translation concepts can bring for enhancing LLMs’ performance in generating professional level translations.

3 Research Design

Building on the effectiveness of introducing contextual and domain-specific information as demonstrated in the literature, this paper investigates prompt design in light of two conceptual tools rooted in translation research: first, “translation brief” as featured in the functionalist approach to translation; second, the “author-translator” dynamic given the persona-matching feature of ChatGPT.

3.1 Prompt design

In total, four prompts were tested in this pilot study, including one basic prompt functioning as a baseline for comparison, and three other prompts featuring three keywords in the scholarship of translation studies: translation brief, author, and translator.

For the basic prompt, because the aim is to evaluate the translation performance of ChatGPT in a professional setting, information included is: 1) a translation command, 2) the target language, and 3) the purpose for professional use, as one would set out in a translation commission. This information was also included in the three other prompts.

For the translation brief prompt, factors including intended text functions, addressees, time and place of text reception, the medium, and the motive (Munday et al., 2022) were included.

For the author-translator dynamic embedded in

the source-target dichotomy, discussions on these two roles and their implications for translation studies have been well documented in the trajectory of translation research. Assigning a persona to ChatGPT is a key feature of the GPT models, and this provides the possibility of incorporating this pair of keywords into prompt design.

Furthermore, the temperature is set at 0.5 for each prompt to constrain the degree of creativity that ChatGPT can potentially exhibit, mimicking the freedom that translators can potentially take in translating articles of this genre in real-life scenarios.

An overview of the four prompts is presented in Table 1.

3.2 Text generation

The source text (ST) selected for the study is a popular scientific article published on the website of the *Discover Magazine* in December 2021.² This genre is chosen for its dual emphasis on maintaining rigorous scientific accuracy and employing a nuanced narrative style, which requires authors and translators to communicate scientific knowledge in a manner that is both accessible and engaging to their respective audiences. The article, titled “A Major Time Travel Perk May Be Technically Impossible”, was written by Cody Cottier, a professional popular science writer. Drawing on a publication of researchers based at the University of Queensland in Australia, the popular scientific article provides accessible and engaging information about time travel for an English-speaking audience interested in but not necessarily have specialized knowledge of this topic.

The selection criteria for the ST are influenced by multiple factors: first, the May 2023 version of ChatGPT utilized in this research has a knowledge cut-off date of September 2021; second, its token capacity (i.e., how many texts it can handle in a single input) is limited; third, the ST should be a professional text; and fourth, a published translation which can serve as a reference document for automatic quality evaluation should be available. To satisfy these basic requirements, the ST is manually checked against the lexical updates on the Oxford English Dictionary website³ to ensure it does not contain any neologisms coined after September 2021. Also, the length of the ST (1253

²<https://www.discovermagazine.com/the-sciences/a-major-time-travel-perk-may-be-technically-impossible>

³<https://www.oed.com/information/updates>

Prompts	Content
Basic	Please translate the following text from English into Chinese Mandarin. The translation is intended for professional use. Top_p=0.5
TransBrief	Please translate the following text from English to Chinese Mandarin. The paragraph is taken from a popular scientific article published in <i>Discover Magazine</i> . The translated version will be published on the <i>Scientific American</i> website in 2023 for professional use. The author of the original text is a well-known science writer, and the target audience for the translation consists of educated individuals interested in popular science. The original text aims to communicate recent research in mathematics that explores the fundamental principles of time travel. Top_p=0.5
Author	You are a professional popular science author. Please translate the following text from English into Chinese Mandarin. The translation is intended for professional use. Top_p=0.5
Translator	You are a professional popular science translator. Please translate the following text from English into Chinese Mandarin. The translation is intended for professional use. Top_p=0.5

Table 1: Prompt overview

words) is manageable for ChatGPT. The authoritative status of *Discovery* in popular science journalism and the availability of a published Chinese translation by *Huanqiuqixue* – a renowned popular science magazine in China – further make the ST a suitable choice.

The model used in the experiment is GPT-4, accessed via the user interface of ChatGPT. Compared to GPT-3.5, this model has demonstrated superior performance in machine translation (Jiao et al., 2023; Wang et al., 2023). All translation outputs were generated by the 24 May 2023 version of ChatGPT. Markdown language was used in the ST to help ChatGPT differentiate headings from main texts and infer the structure of the ST based on the text formatting. Delimiters were used to define the beginning and the end of the ST. Since ChatGPT cannot generate a complete translation in a single response, the prompt “go on” was used to resume the translation command. To assess the consistency of translation outputs generated by the prompts, each prompt was tested three times using a sample sentence from the ST. The outputs were then manually examined by the author for consistency, with a rating scale ranging from 0 to 3, where 0 denotes “Professionally Unusable”, 1 denotes “Professionally Usable with Major Modification”, 2 denotes “Professionally Usable with Minor Modification” and 3 denotes “Professionally Usable”. All four prompts consistently produced similar translations based on the rating. The fourth output from each prompt was selected as the sample for the analysis.

The translation published in *Huanqiuqixue* was labeled as TT1, and four machine translations were labeled as TT2 (Basic), TT3 (TransBrief), TT4 (Author) and TT5 (Translator), where TT stands

for Target Text. The summary of the word count of Chinese characters in each TT (mean \approx 2430, standard deviation \approx 88) is presented in Table 2 below, offering an idea about the size of the translations.

Translation	Word Count
TT1	2602
TT2	2379
TT3	2374
TT4	2369
TT5	2424

Table 2: Summary of the word count of Chinese characters in each translation

3.3 Quality evaluation

Both automatic and human evaluations were conducted to assess the quality of the translation outputs. Two quality evaluation metrics were adopted in this study: BLEU (Papineni et al., 2002) and COMET-22 (Rei et al., 2022). COMET-22 was chosen for its outstanding performance in WMT22 Metrics Shared Task and availability (Freitag et al., 2022). Although BLEU has been criticized heavily for its reliability, it has been chosen as a reference to triangulate results generated by COMET-22 and human evaluations.

To prepare the ST and TTs for automatic evaluation, SDL Trados Studio 2022 was used to align the source and target segments. In total, 66 aligned segments were generated for each ST-TT pair. These aligned texts were then converted into plain text files for BLEU and compiled in an Excel workbook for COMET-22. For BLEU, the text files were processed through the user interface developed by Tilde. For COMET-22

(wmt-comet-da⁴), the metric was run in Python to generate results.

Human evaluations were conducted for qualitative analysis. Four evaluators contributed; all of them are university lecturers based in the UK, who have extensive theoretical and practical knowledge of English-Chinese translation. The evaluators were invited to grade all five TTs (four machine translation outputs and one human translation), without knowledge of which ones were machine-generated translations. Ethical approval was granted by the Humanities and Social Sciences Ethics Committee of Swansea University, before the collection of evaluations (research ethics approval number: 2 2023 6610 5739). Each evaluator was provided with an information sheet and a consent form before taking part in the evaluation.

The grading form designed for human evaluation is different from the metrics typically used in the development of machine translation systems, such as those outlined by Freitag et al. (2022). Instead, it was designed from a translation studies perspective to encourage evaluators to assess the translations on a textual level, following a “top-down approach” (Han, 2020) to obtain a relative ranking of the TTs.

Furthermore, to capture individualized responses regarding the strengths and weaknesses of translations, fixed rubrics containing guided scales were intentionally omitted. This decision stems from the understanding that translation is more than technical transfer of information and that evaluators are not only experienced translation assessors but also readers within this context. Traditional evaluation scales often focus on aspects such as “accuracy” and “adequacy” to ensure replicability and other concerns in machine translation quality assessments. However, such criteria can oversimplify the nuanced nature of translation as a social activity.

Discussions on good versus bad translations are not the primary concern in translation studies; rather, since the cultural turn in the 1990s, translation has been discussed as a socio-historical phenomenon. This viewpoint allows individual interpretations of a ST to be manifested through the medium of translation, which can influence social narratives in another language or culture. This is also true for popular scientific articles embed-

ded with tactical narratives. Traditional criteria reduce the complex social dynamics of translation to mere encoding and decoding of static information, which does not reflect how audiences engage with translated works in real-life scenarios.

Without an evaluation scale that comprehensively considers reader reception, the method adopted in this study allows evaluators the freedom to express their opinions without much interference. This approach provides a more accurate reflection of the real world reception of translations. Admittedly, this might not be the case for some domains, and it would be beneficial to have a reader oriented scale to use, especially at this point of AI development, but it is beyond the scope of the current project.

Based on semantic and structural information embedded in the ST, it was divided into ten segments to create a reading flow for evaluators that resembles the natural reading habits of humans, rather than soliciting evaluations for the sake of evaluation. The source and target segments were aligned in ten blocks in the grading form for easier comparison. Numerical grading boxes (based on a scale of one to ten, with one being the worst and ten the best) and optional free text boxes were provided for each segment. An overall rating block was also included at the end of the grading form. Figure 1 provides a glimpse into the grading form.

In total, each evaluator recorded eleven grades for each TT. For segment grades, the averages were taken for each segment in order to obtain the relative ranking, detailed information can be found in section 4.2.

4 Findings and Discussion

Results from the automatic evaluation metrics and human grading forms provide complementary insights into the quality of the generated TTs, indicating the efficacy of each prompt. This section starts with the results of the two automatic metrics, before delving into human evaluation results.

4.1 Machine evaluation

BLEU and COMET-22 provide scores at both segment and whole text levels. Therefore, each TT yields 67 data points (66 segment scores and one overall score). Table 3 presents the overall scores for the four AI-generated TTs in BLEU and COMET-22, with the rankings shown as superscripts.

⁴<https://huggingface.co/Unbabel/wmt22-comet-da>

Segment 1					
ST	TT1	TT2	TT3	TT4	TT5
A Major Time Travel Perk May Be Technically Impossible	即使能够回到过去，你也改变不了现在：时间旅行的物理学	一项主要的时间旅行优势可能在技术上不可能实现	时间旅行的主要优点可能在技术上不可能实现	时间旅行的主要优势可能在技术上是不可实现的	时间旅行的一大好处可能在技术上无法实现
Grades					
Comments (optional)					

Segment 2					
ST	TT1	TT2	TT3	TT4	TT5
We could, in theory, go back in time. But no matter what we do, the past will likely always lead to the same future.	物理理论允许时间旅行的存在，但无论我们如何改变过去，事情总是会通向同样的未来。	理论上，我们可以回到过去。但无论我们做什么，过去可能总会导致同样的未来。	理论上，我们可以回到过去。但无论我们做什么，过去可能总会导向相同的未来。	理论上，我们可以回到过去。但无论我们做什么，过去可能总会引导到相同的未来。	理论上，我们可以回到过去。但无论我们做什么，过去可能总会引导我们走向相同的未来。
Grades					
Comments (optional)					

[...]

Overall	TT1	TT2	TT3	TT4	TT5
Grade (1 to 10)					
Comments (optional)					

Figure 1: Human evaluation grading form - an example

Metric	TT2	TT3	TT4	TT5
BLEU	4.032 ²	3.314 ⁴	3.93 ³	7.891 ¹
COMET	0.82333 ³	0.82244 ⁴	0.82472 ²	0.82961 ¹

Table 3: BLEU and COMET overall scores and rankings

In both metrics, TT5 (translator) achieved better performance than the three other TTs, and TT3 (translation brief) was ranked the lowest quality. The rankings of the four TTs in BLEU and COMET, however, are different with regard to TT2 and TT4, as shown in Table 3. In general, TT5 (translator) achieved the highest rank across the two metrics, with TT2 (basic) and TT4 (author) following behind. TT3 (translation brief), however, hit the less optimal ground.

Additionally, the differences of the segment scores were tested between TT2 (basic) and TT3 (translation brief), TT2 (basic) and TT5 (translator), and TT4 (author) and TT5 (translator). Wilcoxon matched-pairs signed-ranks tests were employed due to the non-normal distribution of data. Statistical analyses were conducted in Python using the pandas (McKinney and others, 2010) and scipy.stats (Virtanen et al., 2020) packages.

Results show that none of the differences are statistically significant. In BLEU, for the translation brief prompt, the overall score for TT2 (4.03) is higher than TT3 (3.31) by approximately 21.75%. However, the difference, based on the segment scores, is not statistically significant ($p = 0.126$, effect size = 1.21). For the persona group, the over-

all score for TT5 (7.89) is higher than TT4 (3.9) by approximately 102.3%. Yet, the difference at a segment level is also not statistically significant ($p = 0.785$, effect size = 4.06). For the COMET-22 segment scores, results are also insignificant: for TT2 and TT3, the p-value is 0.7853 (effect size = 13.30) and for TT2 and TT5, the p-value is 0.190 (effect size = 0.618). For TT4 and TT5, the p-value is 0.2501 (effect size = 12.73).

These statistically insignificant results could be attributed to the fact that both BLEU and COMET-22 were not initially designed to evaluate the effectiveness of individual prompts within a system. Another potential explanation is that the published translation may not be a suitable reference document for these automatic metrics: even though the omissions and relocations of information in the published translation could potentially enhance its overall communicative effect, this type of translation behavior does not align with the algorithms embedded in BLEU or COM-ET-22. Equally, it could also be the case that the information typically provided in translation briefs does not assist ChatGPT in producing better translations in the same way that it assists human translators. To have a better insight into these issues, the following section reports on human evaluation results.

4.2 Human evaluation

At a document level, the overall grades given by the evaluators and the standard deviations are listed in Table 4 below. No statistical tests were conducted to assess the significance of differences

due to the small number of data points generated in this set of evaluations.

TT No.	Reviewer				Avg	Rank
	1	2	3	4		
TT1	7	9	9	5	7.5	1
TT2	5	4	4	5	4.5	4
TT3	4	4	4	6	4.5	4
TT4	4	6	6	6	5.5	3
TT5	5	6	6	6	5.75	2

Table 4: Human evaluation: Overall scores and rankings

TT1, the published version, received the highest ranking on average. Interestingly, among the four machine translations, human evaluation results also show a preference for TT5 (translator) over the three other prompts. The rankings of TT4 and TT5 also indicate that assigning a persona to ChatGPT tends to enable it to produce a better translation, compared to the translations produced with the basic and the translation brief prompts.

Whilst the overall grades of TT2 and TT3 are identical, the average grades of individual segments reveal a difference between the two. At the segment level, the ten segments add up to a total score of 100. Given that the evaluators for the TTs are the same, taking the average of the segment scores helps to cancel out the individual preferences of each evaluator as a result of maintaining the relative ranking of each translation, based on the assumption that all evaluators are consistent within their own scoring schemes.

Table 5 shows the sums and averages of segment scores for each TTs below. As can be seen in Table 5, the performance of TT5 is the highest among the four prompted outputs, followed by TT2, TT4, and TT3, and these data are in line with the overall scores for the TTs in the automatic metrics.

TT No.	Reviewer				Avg	Rank
	1	2	3	4		
TT1	66	72	89	53	70.25	1
TT2	47	54	52	57	52.50	3
TT3	41	48	58	59	51.50	5
TT4	38	56	57	58	52.25	4
TT5	46	54	57	60	54.25	2

Table 5: Human evaluation: Accumulated sums of segment scores and rankings

Moving on to the comments given by the evaluators for the TTs, for the machine translation outputs, three keywords emerged among the issues pointed out by the evaluators: fluency/naturalness, reader-friendliness, and accuracy.

First, comments on the issues of fluency and naturalness suggest problems associated with syntax, collocation, and lack of creativity in rendering expressions that are not commonly seen in Chinese languages. For instance, the verb “lead” in segment [2] “the past will likely always lead to the same future” was translated as 导致 (lead to a result), 导向 (lead to a direction, usually as a noun) and 引导到 (to guide to) by ChatGPT, which were commented by evaluators on lexical choices that “tend to be made at a surface level”.

Second, taking reader experience into consideration, comments were made on the literal translations of source segments by ChatGPT as “may distract or discourage the readers”, “I’m not sure what this is supposed to mean”, “difficult to follow”, and “this [translation segment] is not clear”. The semantic emphasis of Chinese, especially the use of particles to indicate tenses, also tends to be ignored in the machine translations, as an evaluator mentioned.

Third, two inaccurate translations have been identified by evaluators. For instance, there is one omission example identified by evaluators: a piece of information included in brackets in the ST was omitted in TT2, which led to a fluency issue as an evaluator pointed out, quoting “The text reads more fluently when this clause is included as an organic part of the sentence.” Another case in point is related to terminology accuracy in context. Segment 5 in the ST starts with “no one knows whether time travel is physically possible”, and “physically” here was rendered as 物理上 (literally, regarding Physics) in all four ChatGPT translations. As an evaluator notes, this translation “makes sense but is not as accurate and easy to understand as 技术上” (literally, technically), as seen in the human translation.

For the human translation, on the other hand, most comments are related to the issue of accuracy, specifically with regards to the deviation of meaning and omission cases. This issue, as shown in the comments, is mainly related to the creative modifications of the original text made by the human translator. Creativity, in this case, presents itself as a double-edged sword. For instance, the creative translation of the title was highlighted by evaluators, both as strengths and weaknesses from different perspectives. For one evaluator, the human translation of the title was favored by one evaluator, quoting “I think ‘major time travel perk’ is

difficult to render in Chinese [...] Strictly speaking, TT1 did not follow the ST but adopted a more creative solution. I really like this translation. This sounds exactly like the title of an article you'd read in a popular science magazine." Notably, another evaluator also commented on the positive impact that the freedom shown by the human translators in rendering the title, but at the same time, the negative impact was also pointed out: "It is in the style of title to start with; it conforms less closely to the wording of the ST but incorporates an understanding of the whole article. This is something an experienced translator with good Chinese skills would do or would aim for, at least. Nevertheless, this translation apparently suggests the main purpose of the article is to introduce the physics of time-travel, which is slightly off target." Similarly, in another segment, the translation of a subheading "Time Without Beginning" as 没有起点的故事 (literally, a story without beginning), was pointed out by one evaluator as inaccurate, due to the mis-translation of "time" as "story".

4.3 Summary of Findings

Overall, based on automatic evaluation metrics and human evaluation scores, the rankings of the TTs show that the basic prompt led to better performance of ChatGPT in translation than the prompt including information typical of a translation brief. For the employment of personas to guide ChatGPT, assigning the role of a translator is more effective than the basic prompt and assigning the role of an author, and it has actually led to the best performance among the four prompts tested. For human evaluation comments, it is shown that while the main issues with ChatGPT-generated translations rest on the issues of fluency and naturalness, the comments for the published translation focus mainly on accuracy, mostly resulting from the creativity and stylistic choices shown in the text.

These findings suggest that providing the information contained in a typical translation brief used in human-to-human communication for translation commissions does not necessarily lead to a better performance of ChatGPT in completing translation commands, and that assigning ChatGPT with the role of a translator appears to have a better result than assigning the role of an author or just using a basic prompt.

5 Conclusion

This study explores the efficacy of integrating concepts developed in translation studies into prompting ChatGPT for translation tasks. By evaluating the outputs generated by ChatGPT under four different prompts, it seeks to provide insights into the effectiveness of giving a translation brief to ChatGPT and assigning ChatGPT the personas of an author and a translator. Findings show that assigning the persona as a translator allowed ChatGPT to achieve the best performance among the four prompts, and that the translation generated by ChatGPT using the translation brief prompt received the lowest ranking. This indicates that the classical settings of a translation brief, aiming at human-only workflow, might not work as well as one would expect in a human-machine workflow. However, it would be necessary to revisit the conceptual tools developed in translation studies, considering the development of translation technology and the changing landscape in the industry, so as to further consolidate the relevancy and credibility of translation studies as a discipline. Similarly, training GPT models using aligned source and target texts, paired with translation briefs, and exploring other concepts developed in translation studies could be potentially beneficial.

There are some limitations of the current research. For instance, when testing the consistency of the prompts based on the translation outputs generated by ChatGPT, involving multiple raters, and conducting an inter-rater reliability test would be helpful. Additionally, a reader centered human evaluation metrics and interviews with human evaluators would have been a good complement to the information based solely on the textual analysis of evaluators' comments extracted from the grading form. In addition, using document-level quality evaluation metrics might also strengthen the discussion of the results.

As mentioned in the introduction, this research only provides partial insights into the two general research questions, based on the data collected in this experiment. To further develop this line of research, different prompts conveying information about translation concepts could be examined, across various genres, assessed with a human evaluation scale closer to the reality of translation reading by a larger number of human evaluators. This approach would generate more data, allowing for replication and statistical testing to enhance reli-

bility. Additionally, with the development of Generative AI, research into other LLMs for translation purposes could offer valuable comparative insights for both practitioners and researchers in the field.

Thinking forward, as Hendy et al. (2023) rightly note, although GPT models have promising potential in machine translation, their performance remains underexplored compared to commercial machine translation systems. LLMs are developing rapidly as we write. By extending the scope of translation studies from human-to-human communication to human-machine communication, translation researchers can help to co-shape the future of machine translation and theorize the practice of translation in the new era.

Acknowledgements: I wish to express my gratitude to Dr. Caiwen Wang (UCL/University of Westminster), Dr. Min-Hsiu Liao (Heriot-Watt University), Dr. Yu-Kit Cheung (University of Manchester), and Dr. Yunhan Hu (Durham University) for their valuable evaluation of the translations. Additionally, I extend my thanks to Professor Tong King Lee and the anonymous reviewers for providing insightful feedback on the earlier versions of this manuscript.

References

- Freitag, Markus, Ricardo Rei, Nitika Mathur, Chikui Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Gao, Yuan, Ruili Wang, and Feng Hou. 2023. How to design translation prompts for ChatGPT: An empirical study.
- Gu, Wenshi. 2023. Linguistically informed ChatGPT prompts to enhance Japanese-Chinese machine translation: A case study on attributive clauses.
- Han, Chao. 2020. Translation quality assessment: a critical methodological review. *The Translator*, 26(3):257–273.
- Hendy, Amr, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? a comprehensive evaluation.
- Jiao, Wenxiang, Wenxuan Wang, Jen Tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Is chatgpt a good translator? Yes with GPT-4 as the engine.
- Lee, Tong King. 2023. Artificial intelligence and posthumanist translation: ChatGPT versus the translator. *Applied Linguistics Review*.
- McKinney, Wes et al. 2010. Data structures for statistical computing in Python. In *SciPy*, volume 445, pages 51–56.
- Microsoft. 2023. World trend index annual report.
- Munday, Jeremy, Sara Ramos Pinto, and Jacob Blakesley. 2022. *Introducing translation studies: Theories and applications*. Routledge.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Peng, Keqin, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2023. Towards making the most of ChatGPT for machine translation.
- Rei, Ricardo, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Vilar, David, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting PaLM for translation: Assessing strategies and performance.

- Virtanen, Pauli, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. 2020. Scipy 1.0: fundamental algorithms for scientific computing in Python. *Nature methods*, 17(3):261–272.
- Wang, Longyue, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-level machine translation with large language models. In Bouamor, Houda, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore, December. Association for Computational Linguistics.
- Yamada, Masaru. 2024. Optimizing machine translation through prompt engineering: An investigation into ChatGPT’s customizability.
- Zhang, Biao, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study.
- Zhu, Wenhao, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. Multilingual machine translation with large language models: Empirical results and analysis.

Exploring the Correlation between Human and Machine Evaluation of Simultaneous Speech Translation

Xiaoman Wang
University of Leeds
mlxwang@leeds.ac.uk

Claudio Fantinuoli
KUDO
University of Mainz
fantinuoli@uni-mainz.de

Abstract

Assessing the performance of interpreting services is a complex task, given the nuanced nature of spoken language translation, the strategies that interpreters apply, and the diverse expectations of users. The complexity of this task become even more pronounced when automated evaluation methods are applied. This is particularly true because interpreted texts exhibit less linearity between the source and target languages due to the strategies employed by the interpreter.

This study aims to assess the reliability of automatic metrics in evaluating simultaneous interpretations by analyzing their correlation with human evaluations. We focus on a particular feature of interpretation quality, namely translation accuracy or faithfulness. As a benchmark we use human assessments performed by language experts, and evaluate how well sentence embeddings and Large Language Models correlate with them. We quantify semantic similarity between the source and translated texts without relying on a reference translation. The results suggest GPT models, particularly GPT-3.5 with direct prompting, demonstrate the strongest correlation with human judgment in terms of semantic similarity between source and target texts, even when evaluating short textual segments. Additionally, the study reveals that the size of the context window

has a notable impact on this correlation.

1 Introduction

The assessment of interpreting quality is a common practice in both professional interpreting and academic contexts. The results of these evaluations offer valuable insights for a wide range of stakeholders, including interpreter’s clients, users, practitioners, educators, certification bodies, and researchers (Han, 2022).

Assessing quality in interpretation is a complex endeavor. Quality is not only challenging to measure, but it manifests also an “elusive nature” (Becerra et al., 2013, p. 7) making it difficult to define. The notion of quality in fact may vary from one user to another, introducing a substantial degree of subjectivity in determining what constitutes a good translation of speech. Furthermore, the criteria for quality are contingent upon the type of interpretation involved. For instance, in conference interpreting, the emphasis is generally on the quality of the interpreter’s output, encompassing aspects such as content, language, and delivery. In contrast, within community settings like social and healthcare interpreting, interactional competencies and discourse management play a crucial role in determining what quality is (Kalina, 2012).

Traditionally, the assessment of interpreting performances has been carried out manually, a methodology that comes with its own set of pros and cons. On the positive side, human evaluations offer a holistic view of quality by taking into account various facets of the communication process, thereby delivering a more nuanced understanding of interpreting performance (Pöschhacker, 2002; Becerra et al., 2013). Conversely, manual assessment comes with its own set of challenges, including being labor-intensive, time-consuming

and costly (Wu, 2011). Furthermore, the results often have limited generalizability due to either the restricted scope of the data sampled or the inherent complexities associated with evaluating spoken translation.

In light of the limitations, there has been a growing interest to apply automatic metrics to the evaluation of interpreting performances. While traditional statistical metrics like BLEU have shown limited efficacy in capturing translation quality from a user’s perspective, the emergence of semantic vectors and pre-trained, large-scale generative language models has yielded promising results, especially in the domain of written translation (Kocmi and Federmann, 2023). The application of these metrics is gradually extending to the field of spoken translation as well (Han and Lu, 2021b). However, it must be mentioned that orally translated texts possess certain characteristics that might restrict the efficacy of employing metrics designed for written texts. Interpreters, especially in the simultaneous modality, tend to alter the text more extensively than translators, modifying the structure and omitting parts deliberately as a strategy rather than a deficiency. For example, interpreters may omit part of the original when experiencing cognitive overload, when they cannot comprehend the original message, to name just a few (Korpala, 2012). This non-linearity between source and target texts renders the task of automatic evaluation even more challenging.

The adoption of easy accessible and robust automatic evaluation in interpreting offers several potential applications that could benefit a wide range of stakeholders. Firstly, the ability to provide instant feedback to trainees and practitioners would enable them to quickly assess their performance and pinpoint areas for improvement, also in real-time, creating a faster feedback loop that could substantially accelerate autonomous skill development. Secondly, automatic evaluation might aid organizations in consistently and objectively monitoring the quality of their multilingual services. Thirdly, automatic metrics that correlate with human judgments might serve as a useful tool for the continuous evaluation of machine interpretation.

This study addresses two primary questions: First, is there an automatic metric that aligns closely with human judgment and can thus be used to automate the accuracy evaluation of spoken language translation? Second, do these metrics

evaluate human-generated translations, machine-generated translations, or both more effectively?

The rest of the paper is organised as follows. In Section 2 we present an overview of research in the field of automatic evaluation of interpreting performances. In Section 3, we illustrate our research methodology, our data and the experimental design. Section 3.1 describes the dataset created for this task. Section 3.2 describes the process for human evaluation of the translations while Section 3.3 delves on the process followed for the automatic evaluation. Section 4 presents the results. Section 5 introduces some ethical implications. Finally, Section 6 concludes the paper with some discussion and remarks.

2 Related work

The evaluation of translation quality and in particular of accuracy or information fidelity, i.e. the correspondence between source-language and target-language renditions, has traditionally differed between computer science, with its tradition of abundant use of automatic metrics, and translation and interpreting community, with its focus on manual evaluation as perceived by experts and users.

In computer science, evaluation metrics such as Bilingual Evaluation Understudy (BLEU) (Papineni et al., 2002), National Institute of Standards and Technology (NIST) (Dodgington, 2002), Metric for Evaluation of Translation with Explicit Ordering (METEOR) (Banerjee and Lavie, 2005), and Translation Edit Rate (TER) have been foundational in establishing benchmarks for Machine Translation Quality Estimation (MTQE). BLEU and NIST emphasise n-gram precision, with NIST uniquely weighting distinct n-grams. METEOR integrates both recall and precision, while TER quantifies requisite edits for optimal translation alignment. However, recent scholarly discourse have suggested that these metrics, while valuable, may possess intrinsic limitations in encapsulating the multifaceted subtleties and overarching context of linguistic structures (Fernandes et al., 2023). This acknowledgement has precipitated the exploration of advanced, data-driven methodologies for MTQE without references. Neural networks, characterized by their bio-inspired architectures, emerge as a compelling alternative. These computational structures excel in managing voluminous datasets, discerning intricate patterns, and,

crucially, accounting for the inherent complexities associated with linguistic phenomena.

In the field of neural network architectures, the potential of Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs), and the groundbreaking Transformer for semantic similarity computations has been explored. Of these, Transformer-based models like BERT and GPT have gained considerable academic traction due to their outstanding performance across numerous Natural Language Processing (NLP) tasks (Wang et al., 2023; Kocmi and Federmann, 2023; Fernandes et al., 2023; Clark et al., 2020; Xenouleas et al., 2019; Yang et al., 2019; Brown et al., 2020; Hendy et al., 2023; Clark et al., 2019; Vaswani et al., 2017). Central to these architectures is the concept of embeddings: dense vector representations that capture the semantic essence of words or textual segments. Within this high-dimensional space, vectors situated closely denote semantic relatedness. In translation evaluation, embeddings offer a mediating semantic layer, enabling comparisons between source and target linguistic structures. However, the embeddings landscape is complex. Models from the Universal Sentence Encoder Multilingual (USEM) to the Generative Pre-trained Transformer (GPT) produce embeddings with varied purposes. For example, USEM is geared towards retrieving semantically aligned entities, while GPT emphasizes generating context-rich linguistic constructs. These nuances highlight the need for researchers to thoughtfully choose models aligned with their specific research goals.

In interpreting studies, a subdomain of translation studies dedicated to oral translation, the traditional practice has been to assess accuracy manually, with or without references. Many of the evaluation methodologies are derived by written translation. For the assessment based on references, also known as intra-lingua assessment, the notion of "tertium comparationis" stands as a pivotal benchmark within a particular language (Setton and Motta, 2007). Tracing the historical evolution of this methodology, (Carroll, 1966) stands out as a foundational contributor. He experimented with lay people to ascertain accuracy in translations. Building on Carroll's scale, (Tiselius, 2009) refined the process, integrating references to spoken language and interpreting for intra-lingual assessment. These approaches, while receiving considerable acceptance from seasoned interpreters, are

not without limitations. The impact of cognitive-linguistic factors can potentially alter evaluation results (Han and Zhao, 2021). Moreover, intra-lingual evaluations face challenges in adapting to changing contexts and demographics and may lack a universally acknowledged point of reference (Setton and Motta, 2007).

Academic discussions prioritise evaluation methods for gauging accuracy in inter-language interpreting. These methods range from error analysis, as seen in works by (Gerver, 1969) and (Gile, 1995) that identify translation inaccuracies, to propositional analysis, endorsed by researchers like (Mackintosh, 1983) and (Lee, 1999a; Lee, 1999b; Lee, 2002), which examines textual accuracy. However, these methods present challenges in addressing linguistic subtleties and differing interpretations. More recent research emphasizes grading rubrics, tracing back to (Carroll, 1966), which outline performance across competency tiers and are validated in multiple studies (Han, 2016; Han, 2017; Nia and Modarresi, 2019; Wu et al., 2013). Yet, even with proven reliability, this rubric-based evaluation faces hurdles like the development and validation of rubric descriptors and evaluator inconsistencies.

A few studies have explored the effectiveness of various metrics in evaluating the translation quality or interpreting performances. (Chung, 2020), for instance, pinpointed the strong alignment between human evaluations and scores determined by BLEU and METEOR for German-to-Korean translation. Subsequent studies by (Han and Lu, 2021a) and (Lu and Han, 2022) reinforced the merit of these automated tools. Han and Lu (2021) discerned that METEOR's sentence-level evaluations resonated more with human assessments than broader, text-level evaluations. Conversely, Lu and Han's (2022) exploration, fortified with the integration of the BERT model (Devlin et al., 2019), showcased substantial correlations between human and automated evaluations, underlining the potential of these metrics in assessing interpreting performances. A recent study by (Kocmi and Federmann, 2023) employed Large Language Models like GPT to evaluate translation quality across three language pairs, concluding that only models GPT3.5 and above possess the capability for such translation quality assessment.

While initial studies have underscored positive and moderate-to-strong correlations for MT met-

rics such as BLEU and METEOR, to the best of our knowledge no research has been conducted so far on the use of language models for reference-free interpreting assessment. Our study aims to fill this gap.

3 Data and methodology

3.1 Dataset

The dataset used for the study consists of 12 original speeches in English translated into Spanish, each lasting approximately 5 minutes. These videos were curated from a broader selection of real-life contexts, including lectures, business presentations, live tutorials, and political addresses.¹

Although the corpus size is inherently limited, in order to allow high quality human evaluation, the selection of videos was strategically designed to capture a spectrum of speech features. The speeches were distributed equally in terms of gender, with 6 from male speakers and 6 from female speakers. In addition, the accent of the speakers comprises both native and non-native speakers. The nature of the speeches is diversified into three categories: 4 corporate, 4 political, and 4 general presentations. The speeches comprise 3529 tokens.

For the evaluation purpose, each video was simultaneously interpreted in two ways:

- **Translation H:** Human professional interpreters were engaged. Three interpreters, native Spanish speakers, were involved, each responsible for translating four videos. Simultaneous interpreters were required to interpret the entire video (approximately 5 minutes) to preserve the contextual information essential for accurate interpretation. However, only 2 minutes of the videos were randomly selected for the evaluation.
- **Translation M:** Machine interpretation was carried out by the KUDO AI Speech Translator, the only system available for simultaneous translation available at the moment of writing².

All recordings were automatically transcribed, and the transcriptions were post-edited for accuracy by expert linguists proficient in both English

¹The dataset is available under the Creative Commons 4.0 License at https://github.com/renawang26/Information_fidelity

²www.kudo.ai

and Spanish. The goal of this operation was to make sure that the transcripts did not contain errors of transcription. The transcriptions were manually aligned based on semantic units, a critical step due to the absence of formal punctuation commonly found in written texts. These segments are roughly comparable to sentences or smaller paragraphs. The average length of segments is 29.41 tokens.

3.2 Human evaluation

The human evaluation process was guided by the methodology proposed by (Fantinuoli and Prandi, 2021), which uses a Likert scale to assess two key features of interpretation: accuracy (ability of the translation to convey the meaning of the original) and intelligibility (ability of the translation of being understandable). The two dimensions reflect the main criteria at the core of the product-oriented approach to quality evaluation in Interpreting Studies (Tiselius, 2009). For the purposes of this study, however, the focus was exclusively on the feature of informativeness, i.e. accuracy, leaving the assessment of intelligibility and any other potential feature for future research. One of the advantages of this framework lies in its being user-centric and in line with the corpus-based evaluation already established in Interpreting Studies to assess the quality of human interpretation.

The human evaluation was conducted using a diverse group of 18 evaluators. This consisted of 9 professional interpreters and 9 bilingual individuals without any translation or interpreting experience. The goal was to capture a broad and unbiased evaluation of the translations, taking into account both professional expertise and everyday bilingual proficiency. Each evaluator was assigned 4 videos to evaluate. They were informed that the translations were transcriptions of oral simultaneous interpretations.

For each speech, the raters were asked to assess on a six-point Likert scale first the intelligibility of the output (without a comparison with the source speech nor a comparison between the two outputs), then to evaluate the accuracy of the renditions by comparing each one to the source speech.

An important feature of the evaluation process was the anonymity of the translation sources. Evaluators were not informed whether the translations were produced by a human or by a machine. This was a deliberate step to prevent any evaluation

bias, ensuring that the judgment was strictly based on the quality of the translation, irrespective of the producer.

It is important to point out that with a value of 0.0964 the interater agreement is low (“slight agreement” on a Fleiss’ Kappa scale). This aspect showcases the intrinsic complexity of objectively assessing spoken language translation due to different expectations by the evaluators about what constitutes accuracy. This is an insidious limitation of the evaluation of human interpretation (see Han 2022). While the low agreement between multiple raters is expected, it also limits the generalizability of our findings.

3.3 Machine evaluation

Our approach to the machine assessment of spoken language translation is based on the concept of semantic similarity leveraging sentence embeddings and large language model prompting techniques. The rationale behind using embeddings to measure semantic similarity is multifold, as it proffers a host of advantages, such as the provision of a continuous representation (Mikolov et al., 2013; Pennington et al., 2014), incorporation of contextual information (Devlin et al., 2019; Peters et al., 2018), dimensionality reduction (Roweis and Saul, 2000), applicability of transfer learning (Howard and Ruder, 2018; Raffel et al., 2020), multilingual support (Conneau et al., 2017; Wu et al., 2016), interoperability (Ruder et al., 2019; Artetxe et al., 2018), ease of use (Radford et al., 2018), and state-of-the-art performance (Vaswani et al., 2017; Brown et al., 2020) in NLP tasks. The advantage of using sentence embeddings over MTQE models for assessing interpreting performance lies in the ability of embeddings to evaluate without the need for references. This approach contrasts with some MTQE models, which typically depend on references for quality assessments. Furthermore, sentence embeddings are adaptable across a broad spectrum of languages and text genres, offering a versatile solution for evaluating interpreting performance. This flexibility is beneficial across different domains and contexts, whereas MTQE models often necessitate more specific training data to achieve comparable effectiveness. The fundamental operation of this methodology involves vectorising each sentence in both the source and target texts. Essentially, this means mapping each sentence to its corresponding vectors of real num-

bers, thereby projecting them into a shared multi-dimensional space. This conversion of textual data into numerical format empowers the machine to elaborate the semantics of the texts effectively. The subsequent step involves the calculation of cosine similarity, which serves as a measure of the similarity between each language pair.

We employed three neural network models to carry out sentence embedding: all-MiniLM-L6-v2³, Generative Pre-trained Transformer (GPT), and in particular GPT-Ada model⁴, and Universal Sentence Encoder Multilingual⁵ (USEM). By integrating the all-MiniLM-L6-v2 for its efficient, compact design suitable for multi-language applications, alongside GPT-Ada for their advanced generative capabilities and adaptability, and USEM for its extensive language coverage and cross-lingual semantic understanding, this strategy offers a balanced and comprehensive approach. The embeddings obtained were vectorised sentences in both English and Spanish. The sentence embeddings computed with these models were later used to calculate the Cosine Similarity between the source texts and the translations.

In addition to the models for computing word vectors, we tested another method to compute semantic similarity leveraging the prompt functionality of GPT3.5⁶. The Large Language Model was assigned the task of assessing the semantic similarity between pairs of sentences, one in English and the other in Spanish, using a scale ranging from 1 to 5. An example prompt provided was: “Given the two sentences in English and Spanish, rate from 1 to 5 their similarity, where 1 is not similar and 5 very similar.”

3.4 Computing correlations

The human and automatic assessments were put together in an evaluation matrix. Pearson’s correlation coefficients were leveraged to explore the relationships between human and machine evaluations. Specifically, the correlations between human judgments and the cosine similarities derived from the embeddings of segments from the source speech and their corresponding translations

³<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

⁴<https://platform.openai.com/docs/api-reference/embeddings>

⁵<https://tfhub.dev/google/universal-sentence-encoder-multilingual-large/3>

⁶<https://www.openai.com>

(Translation H and Translation M) generated by three models. were calculated. As stated before, the overarching aim was to probe the feasibility of achieving machine-human parity in the evaluation results.

Since the neural network models chosen for this experiment have limited reliably in computing semantic vectors for long texts, we opted to establish a correlation between human and machine evaluations at the segment level. It is essential to emphasize that similarities values obtained from isolated segment pairs have intrinsic limitations since they are not able to consider accuracy across segments. Thus, the quantity of tokens utilized as context for sentence embeddings could potentially affect the model’s contextual comprehension and, subsequently, the precision of semantic similarity assessments.

To take this into consideration, we examined the effect of “window size”, i.e. the number of segments combined into a single one. For this purpose, we computed similarities for window sizes up to five segments. The systematic variation in window size aimed to shed light on how semantic similarity between human and machine evaluations could potentially be influenced by the availability of cross-segment context.

4 Results

To analyse the data, we devised charts from three perspectives, including the comparison of correlation values among evaluation methods, a comparison between Translation H and Translation M, and correlation values based on window size.

In Figure 1 we compare the distribution of correlation values across all machine evaluation methods, namely GPT-3.5, all-MiniLM-L6-v2, GPT-Ada, and UMSE, for both Translation H and Translation M.

The correlation with GPT-3.5 displays the highest median correlation value. The interquartile range (IQR) is also quite narrow, indicating that the correlation values for this method are consistently high and well-aligned with human evaluations. The correlation with all-MiniLM-L6-v2 has a wide IQR, showcasing varied performance. The median value is close to zero, but there are negative outliers, indicating instances where the machine evaluation is inversely related to human judgment. The correlation values with GPT-Ada are relatively consistent, with a narrow IQR. The

median is slightly above 0.3, which indicates moderate alignment with human judgments. UMSE’s performance seems to be close to GPT-Ada with a median slightly above 0.3. The IQR is a bit larger, suggesting a bit more variability in the correlation values.

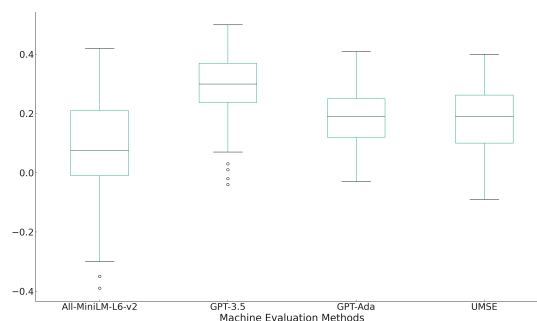


Figure 1: Correlations among machine evaluation methods

For the comparison between Translation H and Translation M in Figure 2, paired bar charts elucidate the average correlation disparities for each machine evaluation method. GPT-3.5’s measurements exhibit robust correlation values with human judgments for both translations, although Translation H marginally outperforms Translation M. For all-MiniLM-L6-v2, the correlation of Translation H gravitates towards zero, whereas Translation M registers a negative value, implying a potential misalignment of the all-MiniLM-L6-v2’s evaluations with human perspectives, predominantly for Translation M. GPT-Ada embeddings yield nearly identical correlation values for both translations, but with Translation H slightly edging out. Intriguingly, UMSE’s embeddings produce a higher correlation value for Translation M compared to Translation H.

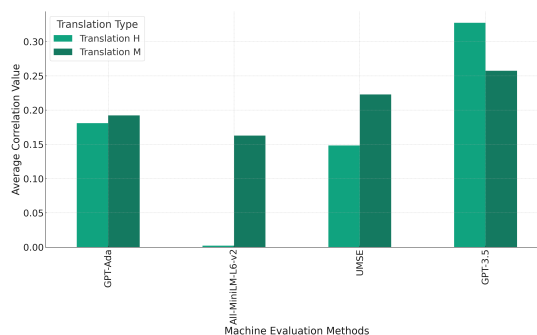


Figure 2: Correlations for Translation H and Translation M

Turning to the third perspective, which examines the shift in correlation values based on window size, line charts in Figure 3 and Figure 4 offer

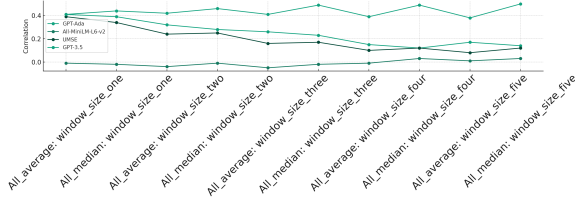


Figure 3: Correlations for Translation H according to window size

insights into this dynamic for each machine evaluation method. For Translation H (Human Translated) in Figure 3, GPT-Ada correlation with human ratings sees a mild fluctuation across window sizes, initially decreasing from window size-1 to size-2, then slightly rising in the following window size 2-5. The all-MiniLM-L6-v2 model, in contrast, exhibits a downward trend, indicating reduced agreement with human evaluations as window size grows. UMSE consistently maintains a stable correlation with human ratings, showing only minor variations across different window sizes. GPT-3.5 presents a distinct pattern; while its correlation initially drops from size-1 to size-2, it surges notably in the subsequent window, outperforming the other models.

In observations for Translation M (Machine Translated) in Figure 4, GPT-Ada begins with a positive correlation with human ratings with size-1, but this declines as the window size expands, hinting at potential metric inconsistencies for broader contexts. The all-MiniLM-L6-v2 model’s correlation, on the other hand, commences positively by size-1 and consistently rises with the window size, pointing to a more aligned evaluation with human judgment for larger translation segments. UMSE’s performance mirrors its evaluation with human translations, maintaining stability across all window sizes and showcasing its consistent metric evaluation. In contrast, GPT-3.5’s correlation fluctuates considerably across window sizes — experiencing a drop, a subsequent rise, and then another decline — indicating a variable level of concurrence with human assessments depending on the window size.

5 Ethical considerations

The adoption of automatic evaluation also raises ethical concerns that warrant careful consideration. One potential issue is the possibility to continuously monitor interpreters, which might infringe

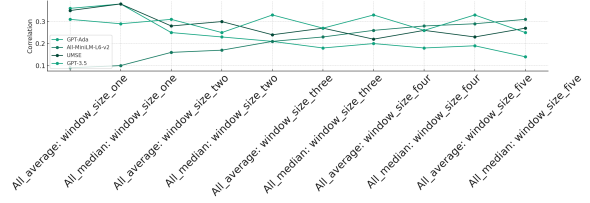


Figure 4: Correlation for Translation M according to window size

on privacy rights and create a sense of constant surveillance, negatively impacting job satisfaction and professional autonomy. Additionally, decisions regarding the employability of individuals could be blindly given to mechanical means, potentially leading to unjust or biased outcomes if the algorithms fail to account for contextual nuances or other essential aspects of human communication. As such, it is crucial for the language industry to carefully weigh the benefits and challenges of automatic evaluation, ensuring that ethical considerations are addressed as advancements in AI technology continue to reshape the landscape.

6 Conclusions

This study aimed to analyze the correlation between automated and human evaluations of translated content. The peculiarity of this experiment is that we focused on a specific form of translation: the simultaneous interpretation of English speeches into Spanish. This mode of translation introduces unique challenges to assessment due to the nonlinear nature of the output (in spoken translation, the differences between the source and target can be more pronounced than in written translation) and varying user expectations regarding interpretation quality. We evaluated both interpretations provided by professional interpreters and those produced by a machine interpretation system. The objective was to develop a metric reflecting interpreting quality in a manner consistent with human judgment.

The direct prompting of GPT-3.5 for quality estimation on a Likert scale exhibits the highest median correlation with human evaluation. This finding establishes GPT-3.5 as the most promising tool among the evaluated methods to gauge translation quality, both for interpretations produced by humans and machines. GPT-3.5 benefits from a larger context, performing better with larger segment windows. This suggests that the model can

capture and evaluate long dependencies more effectively.

Contrary to expectations, GPT-3.5's correlation with human judgment is somewhat stronger for translations produced by professional interpreters than for machine-generated ones. This implies that GPT might be subtly more attuned to the linguistic nuances of human translation, even though it remains adept at evaluating speech translation. One possible explanation for this is that human interpreters often introduce subtle contextual, tonal, and idiomatic adjustments that are more aligned with GPT-3.5's training on diverse data, whereas machine translations might adhere more strictly to equivalences. Looking forward, further research could explore deeper into the characteristics of the datasets used for training such models and their alignment with real-world interpretation tasks.

This study presents several limitations. The observed low interrater agreement suggests potential inconsistencies in human evaluations, possibly affecting correlation values, and generability of the results. Furthermore, the limited scope of the sampled translations might not capture the full range of linguistic complexities inherent to interpretation. Future research should consider evaluations for higher window sizes. In light of GPT-3.5's performance in this study, future research might explore its ability to delineate nuances of typologies of errors rather than merely providing aggregate scores

This study is considered a preliminary attempt to test the feasibility of applying automatic metrics to evaluate inputs from both human and machine interpreters. Before these metrics can be used in production, more research needs to be conducted.

References

- Artetxe, Mikel, Gorka Labaka, and Eneko Agirre. 2018. A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. *arXiv preprint arXiv:1805.06297*.
- Banerjee, Satansjeev and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- Becerra, Olalla García, Macarena Pradas Macias, and Rafael Barranco-Droegge, editors. 2013. *Quality in interpreting. 1*. Number 120 in Interlingua. Comares, Granada.
- Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, and Amanda Askell. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Carroll, John B. 1966. An experiment in evaluating the quality of translations. *Mech. Transl. Comput. Linguistics*, 9(3-4):55–66.
- Chung, H. Y. 2020. Automatic evaluation of human translation: BLEU vs. METEOR. *Lebende Sprachen*, 65(1):181–205.
- Clark, Kevin, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does bert look at? an analysis of bert's attention. *arXiv preprint arXiv:1906.04341*.
- Clark, Kevin, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*.
- Conneau, Alexis, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2017. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv:1810.04805 [cs]*, May. arXiv: 1810.04805.
- Doddington, George. 2002. Automatic evaluation of machine translation quality using n-gram co-occurrence statistics. In *Proceedings of the second international conference on Human Language Technology Research*, pages 138–145.
- Fantinuoli, Claudio and Bianca Prandi. 2021. Towards the evaluation of automatic simultaneous speech translation from a communicative perspective. In *Proceedings of the 18th International Conference on Spoken Language Translation (IWSLT 2021)*, pages 245–254, Bangkok, Thailand (online), August. Association for Computational Linguistics.
- Fernandes, Patrick, Daniel Deutsch, Mara Finkelstein, Parker Riley, André F. T. Martins, Graham Neubig, Ankush Garg, Jonathan H. Clark, Markus Freitag, and Orhan Firat. 2023. The Devil is in the Errors: Leveraging Large Language Models for Fine-grained Machine Translation Evaluation, August. arXiv:2308.07286 [cs].
- Gerver, David. 1969. The effects of source language presentation rate on the performance of simultaneous conference interpreters. In *Proceedings of the 2nd Louisville Conference on rate and/or frequency controlled speech*, pages 162–184. Louisville (Kty), University of Louisville.

- Gile, Daniel. 1995. Fidelity assessment in consecutive interpretation: An experiment. *Target. International Journal of Translation Studies*, 7(1):151–164. ISBN: 0924-1884 Publisher: John Benjamins.
- Han, Chao and Xiaolei Lu. 2021a. Can automated machine translation evaluation metrics be used to assess students' interpretation in the language learning classroom? *Computer Assisted Language Learning*, pages 1–24, August.
- Han, Chao and Xiaolei Lu. 2021b. Interpreting quality assessment re-imagined: The synergy between human and machine scoring. *Interpreting and Society*, 1(1):70–90, September.
- Han, Chao and Xiao Zhao. 2021. Accuracy of peer ratings on the quality of spoken-language interpreting. *Assessment & Evaluation in Higher Education*, 46(8):1299–1313, November.
- Han, Chao. 2016. Investigating Score Dependability in English/Chinese Interpreter Certification Performance Testing: A Generalizability Theory Approach. *Language Assessment Quarterly*, 13(3):186–201, July.
- Han, Chao. 2017. Using analytic rating scales to assess English/Chinese bi-directional interpretation: A longitudinal Rasch analysis of scale utility and rater behavior. *Linguistica Antverpiensia, New Series—Themes in Translation Studies*, 16. ISBN: 2295-5739.
- Han, Chao. 2022. Interpreting testing and assessment: A state-of-the-art review. *Language Testing*, 39(1):30–55, January.
- Hendy, Amr, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.
- Howard, Jeremy and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*.
- Kalina, Sylvia. 2012. Quality in interpreting. *John Benjamins Publishing Company*, 3:134–140.
- Kocmi, Tom and Christian Federmann. 2023. Large Language Models Are State-of-the-Art Evaluators of Translation Quality, May. *arXiv:2302.14520 [cs]*.
- Korpala, Pawel. 2012. Omission in simultaneous interpreting as a deliberate act. *Translation Research Projects 4*, pages 103–111.
- Lee, Tae-Hyung. 1999a. Simultaneous listening and speaking in English into Korean simultaneous interpretation. *Meta: journal des traducteurs/Meta: Translators' Journal*, 44(4):560–572. ISBN: 0026-0452 Publisher: Les Presses de l'Université de Montréal.
- Lee, Tae-Hyung. 1999b. Speech proportion and accuracy in simultaneous interpretation from English into Korean. *Meta: journal des traducteurs/Meta: Translators' Journal*, 44(2):260–267. ISBN: 0026-0452 Publisher: Les Presses de l'Université de Montréal.
- Lee, Tae-Hyung. 2002. Ear voice span in English into Korean simultaneous interpretation. *Meta: Journal des traducteurs/Meta: Translators' Journal*, 47(4):596–606. ISBN: 0026-0452 Publisher: Les Presses de l'Université de Montréal.
- Lu, Xiaolei and Chao Han. 2022. Automatic assessment of spoken-language interpreting based on machine-translation evaluation metrics: A multi-scenario exploratory study. *Interpreting*. ISBN: 1384-6647 Publisher: John Benjamins Publishing Company Amsterdam/Philadelphia.
- Mackintosh, J. 1983. *RELAY INTERPRETATION: AN EXPLORATORY STUDY*. University of London. Ph.D. thesis, unpublished MA thesis.
- Mikolov, Tomas, Ilya Sutskever, Kai Chen, Greg S. Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26.
- Nia, Foroogh Khorami and Ghasem Modarresi. 2019. A Rasch-based validation of the evaluation rubric for consecutive interpreting performance. *Sendebare*, 30:221–244. ISBN: 2340-2415.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Pennington, Jeffrey, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Peters, Matthew E., Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. CoRR abs/1802.05365 (2018). *arXiv preprint arXiv:1802.05365*.
- Pöschhacker, Franz. 2002. Quality Assessment in Conference and Community Interpreting. *Meta*, 46(2):410–425, October.
- Radford, Alec, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. Publisher: OpenAI.
- Raffel, Colin, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text

- transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551. ISBN: 1532-4435 Publisher: JMLRORG.
- Roweis, Sam T. and Lawrence K. Saul. 2000. Non-linear dimensionality reduction by locally linear embedding. *science*, 290(5500):2323–2326. ISBN: 1095-9203 Publisher: American Association for the Advancement of Science.
- Ruder, Sebastian, Ivan Vulić, and Anders Søgaard. 2019. A survey of cross-lingual word embedding models. *Journal of Artificial Intelligence Research*, 65:569–631. ISBN: 1076-9757.
- Setton, Robin and Manuela Motta. 2007. Syntacrobatics: Quality and reformulation in simultaneous-with-text. *Interpreting*, 9(2):199–230, January. Publisher: John Benjamins.
- Tiselius, Elisabet. 2009. Revisiting Carroll’s scales. *Testing and assessment in translation and interpreting studies*, pages 95–121. Publisher: John Benjamins Amsterdam.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Wang, Longyue, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-Level Machine Translation with Large Language Models, April. arXiv:2304.02210 [cs].
- Wu, Jessica, M. H. Liu, and Cecilia Liao. 2013. Analytic scoring in interpretation test: Construct validity and the halo effect. *The making of a translator: Multiple perspectives*, 277:292. Publisher: Bookman Taipei.
- Wu, Yonghui, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, and Klaus Macherey. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Wu, Shao-Chuan. 2011. *Assessing simultaneous interpreting: A study on test reliability and examiners’ assessment behavior*. PhD Thesis, Newcastle University.
- Xenouleas, Stratos, Prodromos Malakasiotis, Marianna Apidianaki, and Ion Androutsopoulos. 2019. Sumqe: a bert-based summary quality estimation model. *arXiv preprint arXiv:1909.00578*.
- Yang, Zhilin, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R. Salakhutdinov, and Quoc V. Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32.

MTUncertainty: Assessing the Need for Post-editing of Machine Translation Outputs by Fine-tuning OpenAI LLMs

Serge Gladkoff², Lifeng Han¹, Gleb Erofeev², Irina Sorokina², and Goran Nenadic¹

¹ The University of Manchester, UK

² Logrus Global, Translation & Localization

lifeng.han, g.nenadic @ manchester.ac.uk

gleberof, irina.sorokina, serge.gladkoff @ logrusglobal.com

Abstract

Translation Quality Evaluation (TQE) is an essential step of the modern translation production process. TQE is critical in assessing both machine translation (MT) and human translation (HT) quality without reference translations. The ability to evaluate or even simply estimate the quality of translation automatically may open significant efficiency gains through process optimisation. This work examines whether the state-of-the-art large language models (LLMs) can be used for this uncertainty estimation of MT output quality. We take OpenAI models as an example technology and approach TQE as a binary classification task. On **eight language pairs** including English to Italian, German, French, Japanese, Dutch, Portuguese, Turkish, and Chinese, our experimental results show that fine-tuned **GPT3.5** can demonstrate good performance on translation quality prediction tasks, i.e. *whether the translation needs to be edited*. Another finding is that simply increasing the sizes of LLMs does not lead to apparent better performances on this task by comparing the performance of three different versions of OpenAI models: *Curie*, *Davinci*, and **GPT3.5** with 13B, 175B, and 175B parameters, respectively.

1 Introduction

Most modern translation projects include post-editing (PE) of machine-translation (MT) output (Han and Gladkoff, 2022; Gladkoff et al., 2022). Instead of translating from scratch, the MT+PE

process increases productivity and allows to speed up global content delivery (Gladkoff and Han, 2022; Han et al., 2013). However, in regulated industries and many other scenarios raw MT output is not suitable for final publication due to the inevitable errors caused by the inherently stochastic nature of neural MT (NMT) (Han, 2022a; Freitag et al., 2021; Hong et al., 2024). Hallucinations, incorrect terminology, factual and accuracy errors, small and large, as well as many other types of mistakes are inevitable to varying degrees of extent, and therefore for premium quality publication human revision is required. MT output serves as input for a professional human translator, who reviews and revises the MT proposals to eliminate factual errors and ensure that the quality of translated material conforms to the customer specifications. At the same time even with those languages that are not handled well by MT, there is a significant portion of segments that are not changed after human review. This portion varies from 10% to 70% in some cases¹, and the question arises, “Is it possible to use machine learning (ML) methods to mark these segments and save time for human reviser and make them focus on those segments that need attention instead”? In other words, *Is it possible to capture editing distance patterns from data of prior editing of this material, which already has been made?* This could further speed up the translation process and decrease the costs while preserving the premium quality of the translated product.

This problem is also closely related to the traditional MT quality estimation (QE) shared task that has been held with the Workshop of MT (WMT) series since 2012 (Callison-Burch et al., 2012; Koehn et al., 2022; Zerva et al., 2022; Han et al., 2013; Han, 2022b), where both token-level and segment-level QE were carried out.

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹logrusglobal.com statistics

From practical application and industrial usage, we formulate the problem into a single classification task, i.e. we are trying to solve a classification task to answer if the translated segment (sentence) needs to be edited, or not.

With the development of current large language models (LLMs), we choose OpenAI models as state-of-the-art LLMs to examine their capabilities for this task. In this work, our first experimental investigation is on “**Predictive Data Analytics with AI: assessing the need for post-editing of MT output by fine-tuning OpenAI LLMs**”. We also follow up with an experiment that explores “**if the size of sample or LLM matters in such a task**” by experimenting with three OpenAI models: *curie*, *davinci*, and *gpt3.5*, with parameter sizes varying from 13B to 175B.

The rest of this paper is designed as below. Section 2 introduces related work to ours including MT-QE-related shared task and challenge events, Section 3 presents our methodology design and pilot study using two language pairs, Section 4 extends the experimental investigation with six more language pairs, section 5 discusses experiment on English-Japanese news content with the increasing sizes of training and testing corpus and explores two more OpenAI LLMs with varying model sizes, and Section 6 concludes this paper with future work and research perspectives.

2 Related Work

The Quality Evaluation (QE) of MT output has always been a critical topic for MT development due to its critical role in assessing quality in the process of training. In many cases, evaluation has to be done without seeing the reference translations. In many practical situations, reference translations are not available or even impossible to acquire, i.e. it is not practical to “manufacture” them for evaluation. The earliest QE shared task with the annual WMT conference started in 2012 when word level QE was introduced by (Callison-Burch et al., 2012) to estimate if the translated tokens need to be edited or not, such as deletion, substitution, or keeping it as it is. In the later development of QE, a sentence-level task was introduced to predict the overall segment translation scores, which are to be correlated with human judgement scores, such as using Direct Assessment (Graham et al., 2015). In WMT-2022, a new task on binary sentence-level classification was also introduced to predict

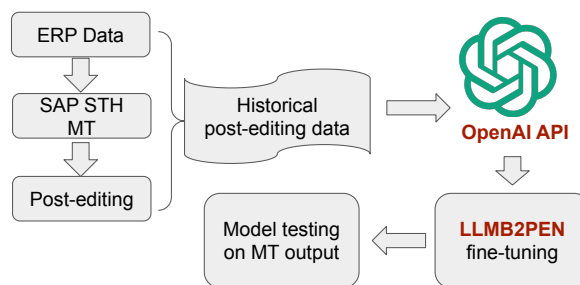


Figure 1: LLMB2PEN Methodology Design on Fine-tuning LLMs for Binary Prediction of Post-editing Need on Translations.

if a translated output has critical errors to be fixed on English-German and Portugueses-English language pairs (Zerva et al., 2022).

The recent methods used for such QE tasks included prompt-based learning using XLM-R by KU X Upstage (Korea University, Korea & Upstage) from Eo et al. (2022), Direct Assessment and MQM features integration into fine-tuning on XLM-R and INFOXLM (Chi et al., 2021) by the Alibaba team (Bao et al., 2022), and incorporating a word-level sentence tagger and explanation extractor on top of the COMET framework by Rei et al. (2022), in addition to historical statistical methods such as support vector machine (SVM), Naive Bayes classifier (NB), and Conditional Random Fields (CRFs) by Han et al. (2013).

However, *to the best of our knowledge*, this work is the first to investigate the OpenAI LLMs with varying sizes on such MT error prediction tasks with positive outcomes.

3 Methodology and Experiments

As shown in the system diagram in Figure 1, we first collect the historical post-editing data from our past projects on eight languages of Enterprise Resource Planning (ERP) content translation on English→German, French, Italian, Japanese, Dutch, Portuguese, Turkish, and Chinese (DE, FR, IT, JA, NL, PT, TR, ZH). This project was completed by using an MT engine to automatically translate the source into the eight languages, followed by post-editing by professional linguists. Two examples of MT and PE in English-Italian and English-German languages as Pilot Experiments are shown in Figure 2 and 3. Regarding MT system selection, since the content was from the ERP domain, we used the SAP STH as our MT

EN-USA	EN-IT MT	MT PE
Have affected personnel in the immediate area of the hazardous work been notified of the work to take place?	Il personale interessato nell'area immediatamente interessata dal lavoro pericoloso è stato informato dei lavori da svolgere?	Il personale interessato dai lavori nell'area vicina all'intervento pericoloso è stato informato sui lavori da svolgere?
Have all Energy Sources been restored per restart procedures and isolations removed?	Sono state ripristinate tutte le fonti energetiche per ogni ripresa e rimozione degli isolamenti?	Sono state ripristinate tutte le sorgenti di energia secondo le procedure di riavvio e sono stati rimossi gli isolamenti?
Have all equipment and safeguards been completely restored?	Tutte le attrezzature e le misure di sicurezza sono state completamente ripristinate?	Attrezzature e salvaguardie interamente ripristinate?
Have all equipment and safeguards been completely restored?	Sono state completamente ripristinate tutte le attrezzature e le salvaguardie?	Attrezzature e salvaguardie interamente ripristinate?
Have all safeguards been completely restored?	Tutte le salvaguardie sono state completamente ripristinate?	Salvaguardie interamente ripristinate?
Have conditions changed?	Le condizioni sono state modificate?	Sono cambiate le condizioni?
Hazardous Energy Control (HEC) Procedural Format	Formato procedurale per controllo energetico pericoloso (HEC)	Formato procedurale per controllo energetico pericoloso (CEP)
Hazardous Energy Control Permit	Permesso di controllo dell'energia pericolosa	Permesso di controllo energia pericolosa (HEC)
Hazardous Energy Control Permit completed as planned	Permesso di controllo dell'energia pericolosa completato come previsto	Permesso CEP (HEC) completato come pianificato
Hazardous Energy Control Permit Template	Modello di permesso di controllo dell'energia pericolosa	Modello di permesso di controllo energia pericolosa (HEC)
Hazardous operations within 35ft (11m) of work area are shut down	Le operazioni pericolose all'interno di 35 ft (11 m) dell'area di lavoro vengono chiuse	Le operazioni pericolose entro 35 ft (11 m) dell'area di lavoro vengono chiuse

Figure 2: EN-IT Examples on MT and Post-Editing

EN-USA	EN-DE MT	MT PE
Chemical/Liquid Protective Boots	Chemische/flüssige Schutzboote	Schutzstiefel für Chemikalien und Flüssigkeiten
Child Location	Unterlokation	Untergeordneter Standort
Child Permit	Kindgenehmigung	Untergeordnete Genehmigung
Child State	Untergeordneter Bundesstaat	Untergeordneter Zustand
Child Workflow	Untergeordneter Workflow	Untergeordneter Arbeitsablauf
Child Workflow State	Untergeordneter Workflow-Status	Untergeordneter Arbeitsablaufstatus
Chilled Water - specify Pressure	Kühlwasser - Geben Sie den Druck an.	Kühlwasser - Druck angeben
Choose Workflows to Import	Zu importierende Workflows auswählen	Zu importierende Arbeitsabläufe auswählen
The file size exceeds the limit allowed	Die Dateigröße überschreitet das zulässige Limit	Die Dateigröße überschreitet die zulässige Grenze
Qualified Person returned service to full operation and notified all affected workers?	Qualifizierte Person hat den Dienst in Betrieb genommen und alle betroffenen Mitarbeiter benachrichtigt?	Hat die qualifizierte Person den Betrieb wieder vollständig hergestellt und alle betroffenen Mitarbeiter benachrichtigt?
Qualified persons are First Aid/CPR Certified	Qualifizierte Personen sind Erste-Hilfe-Zertifizierung/CPR-zertifiziert	Qualifizierte Personen verfügen über Erste-Hilfe-Zertifizierung/CPR-Zertifizierung
Record results every 30 minutes for length of work.	Erfassen Sie die Ergebnisse alle 30 Minuten für die Dauer der Arbeit.	Ergebnisse alle 30 Minuten für die Dauer der Arbeit erfassen.
Record Worst Case Reading of any meter at any level in the space.	Erfassen Sie den Worst Case Reading eines beliebigen Zählers auf einer beliebigen Ebene des Bereichs.	Aufzeichnung des Worst-Case-Messwerts eines beliebigen Zählers auf einer beliebigen Ebene im Raum.
The symbol displayed at the top left of the application, used as the home button.	Das oben links in der Anwendung angezeigte Symbol, das als Home-Drucktaste verwendet wird.	Das oben links in der Anwendung angezeigte Symbol, das als Startseiten-Schaltfläche verwendet wird.

Figure 3: EN-DE Examples on MT and Post-Editing

engine. ²

With this data from a real-world translation project, we used API to fine-tune the OpenAI *curie* model for our classification task. The input is the triple set (English source, MT outputs, post-edited "gold standard") we prepared in Phase 1. The goal of this step is to optimise the weights of the model parameters for our classification task. The custom fine-tuned model produced as a result of LLMB2PEN (LLM for Binary Prediction of Post-editing Need) method is created in our private space on the OpenAI account.

We *did not* apply "prompt engineering" for this task by doing zero-shot, one-shot, or few-shot training; we *did a full-scale fine-tuning* of OpenAI LLMs via API. It is important to note that we did not simply train the LLM for edit distance either; instead, the model was trained to learn whether the strings were edited or not taking into account the full content of the string and the entire context of the training data. One of the reasons that we did not use prompting is that "Prompt Engineering" of ChatGPT-3 is limited by 3,000 tokens, and

²<https://www.sap.com/> SAP is an enterprise resource planning, automation and business software company.

with ChatGPT-4 the context has been increased to 25,000 tokens, but still very significant limitation remains. OpenAI documentation states that 100 tokens = 75 words, meaning that the average sentence is 20 tokens, therefore 3000 tokens is only 150 sentences, or 75 translation units of bilingual text, or 50 segment triples of source, target and reference. The context of 25,000 sentences is only about 150 segment triples.

Also, fine-tuning is a deeper process of adjusting the model's weights, and not just an in-context learning. That's why we chose fine-tuning method, which is not constrained by such limitations.

For our classification experiment we took about 4000 lines of bilingual data in triples of source, target, and reference, and split it into train (large) and test (smaller) sets with a ratio of 9:1.

There were no specific selection criteria for the data because we took the entire project dataset after project completion. (Please, note that since we used the entire data from the actual project, and split the data set as 9:1, the sizes of test sets are not round and slightly different for different languages.)

We also combined source sentences in groups of length, so that the test data set has the same distribution of sentences by their length as the training dataset.

Since the average sentence size is about 17 words, the training dataset contained about 35000 words of source data, 35000 words of MT output, and 35000 words of post-edited human reference.

It is also important to note what the model learns in this case - in such an experiment it learns not to translate, but to spot MT translation errors that were made by the specific MT engine in a specific language pair on particular content.

3.1 Outputs on EN-DE/IT

As a first step, we trained the *curie* LLM model using our data for two language pairs: English-Italian and English-German. To illustrate the results of prediction with our LLMB2PEN method, we draw the confusion matrix for both language pairs in Figures 4 and 5.

In the Confusion Matrix, from the top left corner in a clockwise direction, the 1st quadrant means True Negative (TN): segment is predicted as not requiring editing and it does not indeed require post-editing. The 2nd quadrant is False Positive (FP): segments which are predicted as requiring editing, but in reality, they do not, that is **FP** means that the segment is correct but wrongly flagged for post-editing. The 3rd quadrant is True Positive (TP) - reflecting the situation when a segment is correctly flagged as requiring post-editing. The fourth quadrant is False Negative (FN): segment is predicted as correct, while in reality, it does require post-editing. So the first and third are successful classifications, and the other two are incorrect classifications.

It is worth mentioning that if the segment is incorrectly predicted as requiring post-editing, this only leads to a small increase in post-editing cost, while False Negative predictions represent the consumer’s risk of seeing substandard segments as not corrected in the final product. So in the context of our task, we are much more concerned with the share of False Negatives in the test classification dataset.

In the Italian situation shown in Figure 4, you can see that the model predicts correctly that many more translated sentences need to be edited (TP=503) than sentences that do not need to be edited (TN=191). In incorrectly predicted cate-

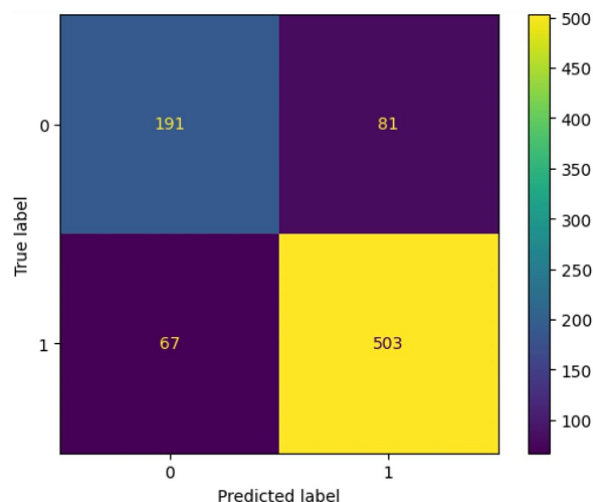


Figure 4: EN-IT Confusion Matrix of LLMB2PEN, *curie* model: Clockwise from top-left corner (TN, FP, TP, FN)

gories, 67 sentences need to be edited but predicted as good, and 81 translated sentences do not need to be edited, but the prediction says they have to be reviewed.

In the English-German set from Figure 5, the situation is the opposite: there are more translated sentences that do not need to be edited (442) than prescribed for review (256) in the correct predictions. In the wrong prediction categories, such numbers are 90 and 46 respectively.

The prediction **accuracy** of the LLMB2PEN model on our designed task is $(TP+TN)/Total = (503+191)/842 = 82.42\%$ for English-Italian MT, and $(442+256)/834 = 83.69\%$ for English-German MT. Overall, our LLMB2PEN method shows that the English-German output is clearly better than the English-Italian.

However, if we only count the Type II errors (incorrect prediction that the segments should NOT be edited), then the corresponding error rates will be $67/842 = 8\%$ for Italian and $90/834 = 10\%$ for German.

3.2 Discussion

The first and foremost finding is that the fine-tuned model learned enough information to make a very significant prediction of whether the segment has to be edited or not. It should be noted that such successful classification holds the promise of a viable method to significantly reduce the volume of post-editing efforts and therefore time and costs. There is, however, a problem: while it is OK to

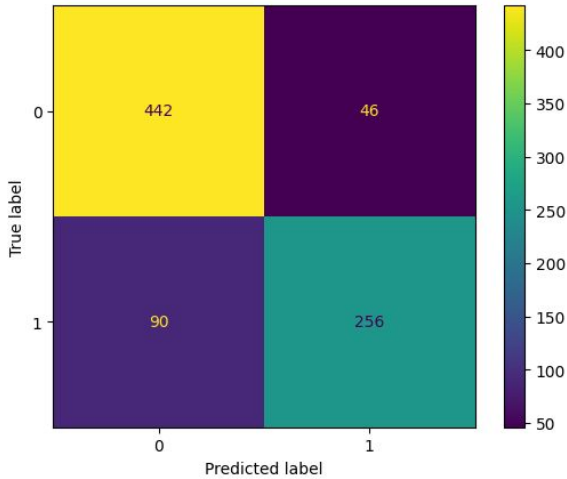


Figure 5: EN-DE Confusion Matrix of LLMB2PEN, *curie* model: Clockwise from top-left corner (TN, FP, TP, FN)

present the editor with segments that are predicted as required for editing, but in reality do not require editing (the fourth quadrant, FP), real consumer risk comes from the segments that have been predicted as not requiring editing and made their way to the final predict, but in reality, they contain errors (the fourth quadrant, FN).

Such segments represent a significant portion of segments predicted as not requiring post-editing: $\text{FN}/(\text{TN}+\text{FN}) = 67/(191+67) = 67/258 = 26\%$ of “leave as is” (let’s call them “LAI”) segments for Italian, and $90/(442+90) = 90/532 = 16.9\%$ for German.

It is possible that for specific language pairs and MT engines the portion of the LAI segments will decrease with the size increase of the training data and further fine-tuning, but it is unlikely to become zero, since with neural models the error rate is never zero.

Two strategies can be considered for implementing such prediction in production:

1. The LAI segments are excluded from the human loop and go into publication unvetted, but not straight away as they advance through the workflow along with all the other segments. In this scenario, the potential error rate ceiling for final content will be $\text{FN}/\text{Total} = \text{FN}/(\text{TP}+\text{FN}+\text{TN}+\text{FP}) = 8\%$ for Italian, i.e. $67/(81 + 67 + 191 + 503) = 81/842$ and $10.8\% = 90 / (90 + 46 + 442 + 256) = 90/834$ for German.

It is not impossible to predict what would be the actual error rate in those 8% and 10.8% segments that will not be reviewed or the severity of errors in them. It is, obviously, the decision of the customer to decide whether this is an acceptable level of consumer risk for their situation (domain, type of content, audience, etc.). Additional risk assessment may be required to be carried out.

The savings on post-editing volume in this scenario would be $(\text{TN}+\text{FN})/\text{Total} = (191+67)/842 = 30.1\%$ for Italian and $(442+90)/834 = 63.8\%$ for German.

2. All LAI segments are marked as “100% MT matches” in a CAT tool. With this approach, translators are requested to review them, but at a lower per-word rate, using the traditional approach which is well familiar to translation providers. In this scenario the reduction of the total time, effort, and cost can be estimated as follows: without this approach, translators working on the Edit Distance Calculation (EDC) model will get lower payment (which can vary from 10% to 40% with different payment models) for not changed segments. In this scenario, translators may be asked to review such LAI segments but paid only a small part of the full rate for the review of such segments.

Simple proportion allows us to calculate the savings in the second scenario: if we take the full payment for all the segments for 100% of post-editing costs, and assume that 10% pay reflects adequate pay for the review of LAI segments that are marked as such, the volume of post-editing decreases 27.6% for Italian and 57.4% for German with zero error rate of the final product (no producer’s or consumer’s risk).

This estimate of a potential economy with a guarantee of zero error rate begs for further research and implementation of this method.

4 Extended Experiments On Six More Language Pairs

We hereby also present extended experimental results using six more language pairs obtained with LMB2PEN method for translation editing distance prediction. These language pairs include English-to-French, Japanese, Dutch, Portuguese, Turkish, and Chinese (EN→FR/JA/NL/PT/TR/ZH), whose

results are listed in Figure 6, 7, 8, 9, 10, and 11 respectively.

From the results presented in the figures, in general, the ratio of correct prediction (TP+TN) is much higher than the one from mis-prediction (FN+FP) across all these language pairs, as for English-Italian and English-German in the pilot studies. On one hand, the following language pairs have more True Positive than True Negative predicted segments **than for English-German/Italian**: English-Japanese, English-Portuguese, and English-Chinese. On the other hand, the rest of the language pairs have more TN than TP: English-French, and English-Dutch, except for English-Turkish which has a comparable number of segments between TP (347) and TN (353) labels. This finding also indicates that such language pairs with a high number of TN labels are still much more challenging for MT system development to produce more correct outputs, i.e., English to French, Dutch, and Turkish. Earlier research findings from Gladkoff et al. (2022) on TQE conclude that 200+ segments can be enough amount of data to reflect the MT system quality.

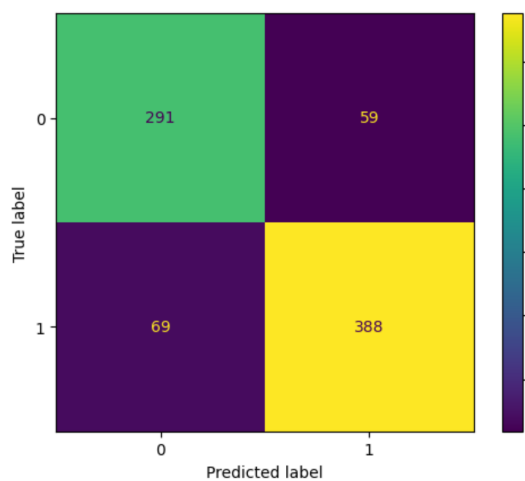


Figure 6: EN-FR Confusion Matrix of LLMB2PEN, *curie* model: Clockwise from top-left corner (TN, FP, TP, FN)

5 Different LLMs on EN-JA News Domain

In the subsequent experiment on data, we used different news items translation corpus from different projects, translated from English to Japanese.

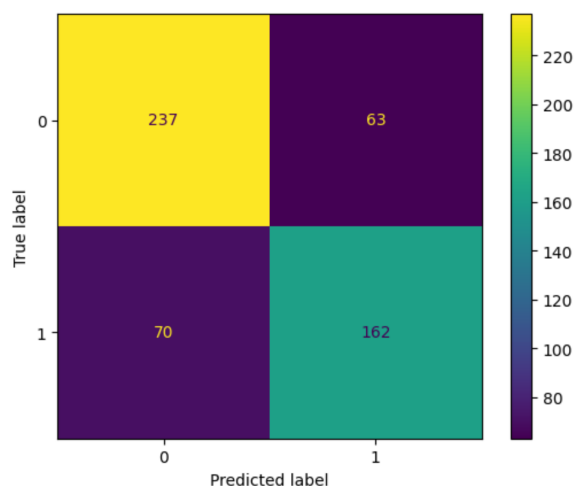


Figure 7: EN-JA Confusion Matrix of LLMB2PEN, *curie* model: Clockwise from left-up corner (TN, FP, TP, FN)

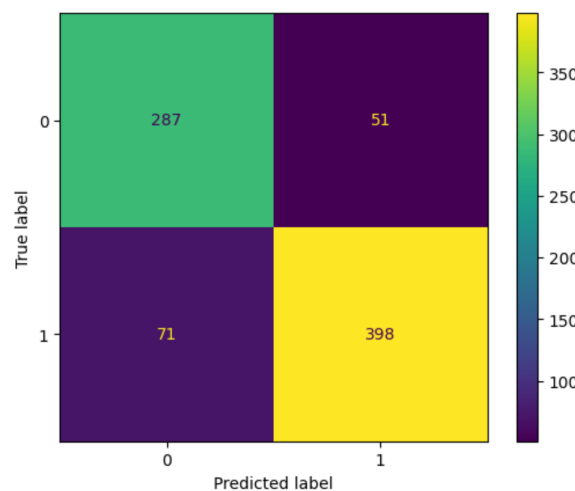


Figure 8: EN-NL Confusion Matrix of LLMB2PEN, *curie* model: Clockwise from top-left corner (TN, FP, TP, FN)

5.1 Using OpenAI GPT3.5turbo

In this experiment, we have repeated experiments of fine-tuning the OpenAI *gpt3.5turbo* model on datasets of different sizes: 2000 pairs, 4000 pairs, and 6000 pairs.

Figure 12 shows the confusion matrix for the training set of 6000 bilingual EN-JA translation pairs in the news domain.

We ran several experiments with varying training set sizes, with results shown in Figure 13.

These results are interesting because although False Positive prediction does not improve with the increase of training set, in the context of the

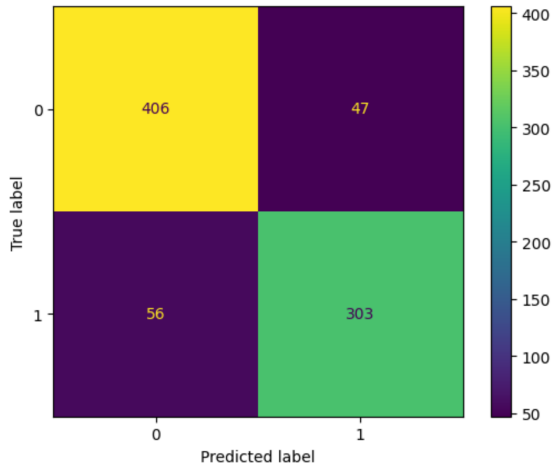


Figure 9: EN-PT Confusion Matrix of LLMB2PEN, *curie* model: Clockwise from left-up corner (TN, FP, TP, FN)

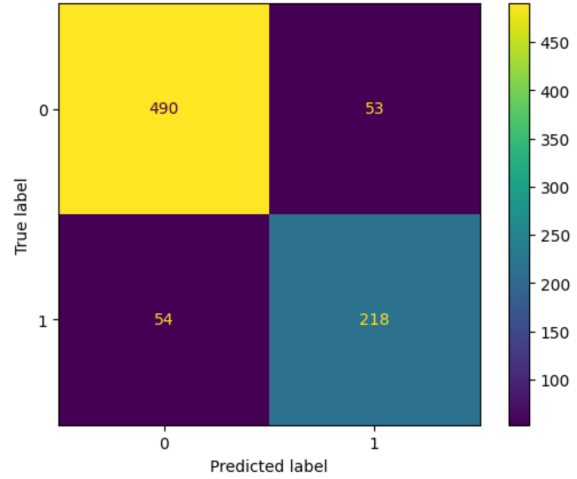


Figure 11: EN-ZH Confusion Matrix of LLMB2PEN, *curie* model: Clockwise from top-left corner (TN, FP, TP, FN)

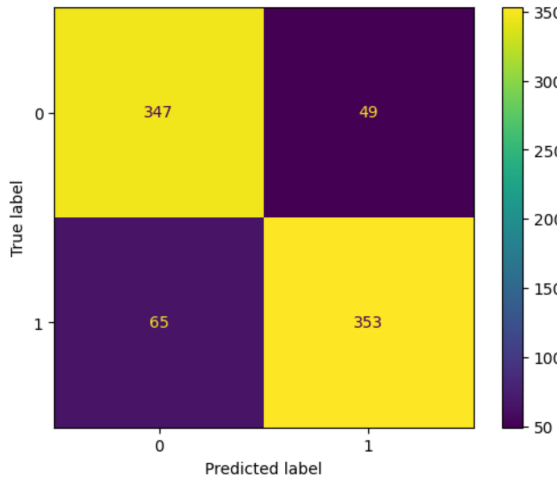


Figure 10: EN-TR Confusion Matrix of LLMB2PEN, *curie* model: Clockwise from top-left corner (TN, FP, TP, FN)

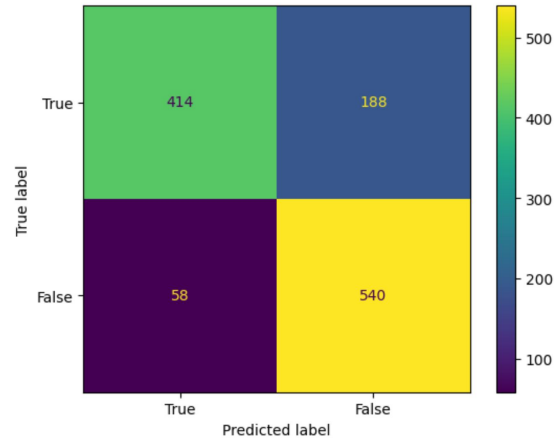


Figure 12: EN-JA news items Confusion Matrix of LLMB2PEN, *gpt3.5turbo* model: Clockwise from top-left corner (TP, FN, TN, FP)

need for post-editing the False Negative category is much more interesting, because we are interested in better prediction of those segments which do NOT require post-editing. And, as we see from the experimental data, the prediction of FN improves from almost 20% to 12%-15% with the increase of training set from 2000 bilingual segments to 6000 bilingual segments.

We, therefore, can recommend the training set in that range, since larger sizes of training set will be more expensive and will take significant time for models with the size of *gpt3.5turbo*.

Train	Test	Total	Legth groups	FN	FN, % total	FP	FP, % total
2000	400	2400	10	79	19.75%	17	4.25%
2000	400	2400	15	74	18.50%	23	5.75%
2000	400	2400	5	74	18.50%	18	4.50%
4000	800	4800	5	162	19.00%	16	2.00%
4000	800	4800	10	102	12.75%	75	9.38%
6000	1200	7200	10	188	15.67%	58	4.83%

Figure 13: EN-JA news items predictions with fine-tuning completed on different training dataset sizes, *gpt3.5turbo* model

5.2 Comparison of performance on different OpenAI models

It was also interesting to see how the extra-large LLMs (xLLMs) from OpenAI, the *davinci* and *gpt3.5turbo* models, perform on the same task in

comparison to *curie* model we used earlier. These three LLMs have parameter sizes around 13B, 175B, and 175B respectively.

So we used the same English-Italian data from our original experiment to compare performance on different models of the same EN-IT dataset.

Figure 14 shows the comparison of these three LLMs regarding their confusion matrix and parameter sets. Surprisingly, their performances on predicting MT errors are very close, i.e. the larger-sized *davinci* model and extra-large sized *gpt3.5turbo* did not demonstrate much improvement on model classification accuracy. Their correct labels (TP+TN) are (694, 699, 706) respectively out of 842 all labels, which results in the accuracy ratios 82.42%, 83.02%, and 83.85%. In comparison to the much smaller *curie* model with 12 layers of Transformer and 768 hidden units, the xLLM *gpt3.5turbo* only achieved 1.43 points (83.85%-82.42%) increase of accuracy score despite using 175 layers of Transformer and 4096 hidden units.

The explanation for this may probably be found because the fine-tuning loss on this classification task drops down very quickly.

Figure 15 shows the fine-tuning loss on the *gpt3.5turbo* model. As can be seen from this graph, only 100 steps are sufficient to bring the loss to almost zero, and then all other steps contribute very little to the classification quality improvement.

As we can see, there is no need to use larger models since results hardly improve as compared with *curie* model.

6 Conclusions and Future Work

In this work, to investigate the LLM's capability of predicting MT output errors, we fine-tuned GPT models via OpenAI API. We formulated the task as a classification challenge using prepared historical post-editing data on English-Italian and English-German for pilot studies. The experimental output using fine-tuned LLMB2PEN demonstrated promising results. We also analysed the possible solutions for addressing the error rates, i.e. whether prediction errors can be ignored and published without the review, or letting them be reviewed by the linguists at a lower rate, and how much saving can be achieved for the client who uses this process, in comparison to 100% post-editing without using LLMB2PEN method.

In the extended experiments, we added six more language pairs including English-to-French, Japanese, Dutch, Portuguese, Turkish, and Chinese, in total resulting in eight, and summarised our findings by classifying the language pairs. We also compared GPT models from different sizes and the experimental results surprisingly show that the larger LLMs (*davinci* and *gpt3.5turbo*) do not improve the accuracy performance of much smaller *curie* model with apparent margins but with much more cost.

In the future, we are going to work on response rate and training times to see whether the model can continue learning as *being fed with more consecutive chunks of data* for the same languages, to implement an ongoing learning of prediction. In addition, we plan to carry out the LLMB2PEN fine-tuning on other language pairs for which we have historical data. We intend to explore to what extent the model is capable of absorbing data for several languages, i.e. one fine-tuned multilingual model serving several language pairs.

To further extend this project, it will also be interesting to explore and check whether the LLMB2PEN method can help to identify human-introduced errors or translationese.

Limitations

In this work, we reported MT QE experiments using eight language data translated from English. The positive results produced from the OpenAI models can be further enhanced by more language pairs, as well as broader domains of the corpus.

The main limitation of the method is non-zero fine-tuning time. The fine-tuning takes about 20 minutes and therefore cannot be made continuous, which has to be done periodically, in batches. This hardly can be overcome, but deployment methods can be applied to quickly replace the older fine-tuned models with the newer ones.

Ethical Statement

This work has no ethical concerns since we did not disclose any identifiable private user data. All experiments were carried out in a secure computing environment.

Acknowledgements

We thank Georg Kirchner, Globalization Technology Manager at Dell Technologies, for the valuable comments on the initial manuscript. LH and

	Curie			DaVinci			ChatGPT 3.5 Turbo		
True Label									
FALSE	0	191	81	201	68	206	63		
TRUE	1	67	503	68	505	80	493		
Predicted Label		0	1						
Total			842		842		842		

Parameter	Curie	DaVinci	GPT-3 (ChatGPT)
Architecture	Transformer	Transformer	Transformer
Layers	12	24	175
Hidden Units	768	1024	4096
Parameters	~13B	~175B	~175B
Training Data	Web Text	Web Text	Web Text
Use Case	Specialized tasks like translation, summarization	General-purpose, wide range of tasks	General-purpose, wide range of tasks
Released	2019	2020	2020

Figure 14: Comparisons of three LLMs on Confusion Matrix and Parameters.

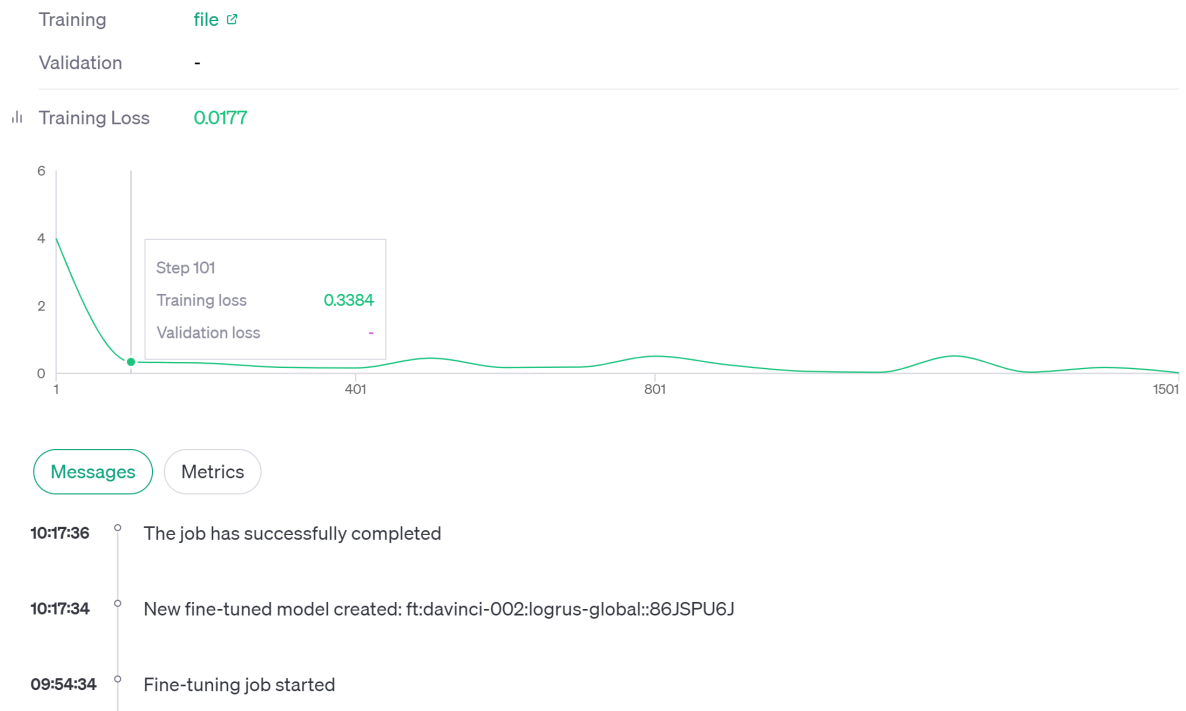


Figure 15: Fine-tuning progress on *gpt3.5turbo* model fine-tuning.

GN are grateful for the support from the grant “Assembling the Data Jigsaw: Powering Robust Research on the Causes, Determinants and Outcomes of MSK Disease”. The project has been funded by the Nuffield Foundation, but the views expressed are those of the authors and not necessarily the Foundation. Visit www.nuffieldfoundation.org. LH and GN are also supported by the grant “Integrating hospital outpatient letters into the

healthcare data space” (EP/V047949/1; funder: UKRI/EPSRC).

References

Keqin Bao, Yu Wan, Dayiheng Liu, Baosong Yang, Wenqiang Lei, Xiangnan He, Derek F. Wong, and Jun Xie. 2022. Alibaba-translate China’s submission for WMT 2022 quality estimation shared task. In *Proceedings of the Seventh Conference on*

- Machine Translation (WMT)*, pages 597–605, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Chris Callison-Burch, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, and Lucia Specia, editors. 2012. *Proceedings of the Seventh Workshop on Statistical Machine Translation*. Association for Computational Linguistics, Montréal, Canada.
- Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao, Heyan Huang, and Ming Zhou. 2021. InfoXLM: An information-theoretic framework for cross-lingual language model pre-training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3576–3588, Online. Association for Computational Linguistics.
- Sugyeong Eo, Chanjun Park, Hyeonseok Moon, Jaehyung Seo, and Heuseok Lim. 2022. KU X upstage’s submission for the WMT22 quality estimation: Critical error detection shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 606–614, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, Errors, and Context: A Large-Scale Study of Human Evaluation for Machine Translation. *arXiv e-prints*, page arXiv:2104.14478.
- Serge Gladkoff and Lifeng Han. 2022. HOPE: A task-oriented and human-centric evaluation framework using professional post-editing towards more effective MT evaluation. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 13–21, Marseille, France. European Language Resources Association.
- Serge Gladkoff, Irina Sorokina, Lifeng Han, and Alexandra Alekseeva. 2022. Measuring uncertainty in translation quality evaluation (TQE). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1454–1461, Marseille, France. European Language Resources Association.
- Yvette Graham, Timothy Baldwin, and Nitika Mathur. 2015. Accurate evaluation of segment-level machine translation metrics. In *NAACL HLT 2015, The 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, Colorado, USA, May 31 - June 5, 2015*, pages 1183–1191.
- Aaron Li-Feng Han, Yi Lu, Derek F. Wong, Lidia S. Chao, Liangye He, and Junwen Xing. 2013. Quality estimation for machine translation using the joint method of evaluation criteria and statistical modeling. In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 365–372, Sofia, Bulgaria. Association for Computational Linguistics.
- Lifeng Han. 2022a. *An investigation into multi-word expressions in machine translation*. Ph.D. thesis, Dublin City University.
- Lifeng Han. 2022b. An overview on machine translation evaluation. *arXiv preprint arXiv:2202.11027*.
- Lifeng Han and Serge Gladkoff. 2022. Meta-evaluation of translation evaluation methods: a systematic up-to-date overview. In *Tutorial at LREC2022*, Marseille, France.
- Kung Yin Hong, Lifeng Han, Riza Batista-Navarro, and Goran Nenadic. 2024. Cantonmt: Cantonese to english nmt platform with fine-tuned models using synthetic back-translation data.
- Philipp Koehn, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors. 2022. *Proceedings of the Seventh Conference on Machine Translation (WMT)*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates (Hybrid).
- Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Chrysoula Zerva, Frédéric Blain, Ricardo Rei, Piyawat Lertvittayakumjorn, José G. C. de Souza, Steffen Eger, Diptesh Kanojia, Duarte Alves, Constantin Orăsan, Marina Fomicheva, André F. T. Martins, and Lucia Specia. 2022. Findings of the WMT 2022 shared task on quality estimation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 69–99, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Translators' perspectives on machine translation uses and impacts in the Swiss Confederation: Navigating technological change in an institutional setting

Paolo Canavese

School of Applied Language and
Intercultural Studies
Dublin City University, Dublin, Ireland
paolo.canavese@dcu.ie

Patrick Cadwell

School of Applied Language and
Intercultural Studies
Dublin City University, Dublin, Ireland
patrick.cadwell@dcu.ie

Abstract

New language technologies are driving major changes in the language services of institutions worldwide, including the Swiss Confederation. Based on a definition of change management as a combination of adaptation measures at both the organisation and individual levels, this study used a survey to gather unprecedented quantitative data on the use and qualitative data on the perceptions of machine translation (MT) by federal in-house translators. The results show that more than half of the respondents use MT regularly and that translators are largely free to use it as they see fit. In terms of perceptions, they mostly anticipate negative evolutions along five dimensions: work processes, translators, translated texts, the future of their language services and job, and the place of translators within their institution and society. Their apprehensions concern MT *per se*, but even more the way it is seen and used within their organisation. However, positive perspectives regarding efficiency gains or usefulness of MT as a translation aid were also discussed. Building on these human factors is key to successful change management. Academic research has a contribution to make, and the coming together of translation and organisation studies offers promising avenues for further research.

1 Introduction

Recent, fast-paced technological developments in the language industry, notably the advent of neural machine translation, are imposing structural and organisational changes in the language services of companies and institutions worldwide. For example, the results of the ELIS 2024 survey indicate that more than 70% of translation departments in national and local agencies and almost all translation departments in international public agencies that participated in the survey have implemented or are implementing MT (ELIS, 2024: 37).¹ Not surprisingly, this wave of change has also hit the Swiss Confederation (Nussbaumer, 2020). Because Switzerland is a multilingual country, the federal institutions rely on a network of language services (LS) within the government and Parliament to communicate in the four national languages – German, French, Italian and Romansh – as well as in English. These LS employ a total of 481 staff members, including translators, heads of service, legal drafters, terminologists, language technology specialists, and trainees.²

Initial, structured, large-scale attempts to integrate MT into the work processes of the Swiss Confederation started in 2019, when 130 DeepL Pro licenses were bought and a specific working group was formed to carry out a test phase. The working group produced an extensive report and a set of recommendations, concluding that MT can be a helpful tool, to be used according to the

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹ Although the actual use is lower, at 34% and 38% respectively (ibid.: 38).

² These statistics were provided by Franco Fomasi (Federal Chancellery) in a private communication with the first author. They refer to the year 2022.

principle: “What you need for yourself, you can machine translate, what you need for others, you better give to a professional!” (Arbeitsgruppe Maschinelle Übersetzung, 2019; see also the relevant press release: Federal Council, 2019). Since then, all staff within the Swiss Confederation, including the LS, has had access to DeepL Pro. At the same time, a Centre of Expertise for Language Technologies (CoELT)³ was set up in November 2020. In addition to providing support and training on CAT tools, MT, and other language technologies, it “keeps track of market developments and, with user involvement, initiates evaluation and procurement projects for appropriate technologies” (ibid.) in order to find the solutions that best meet the needs of the federal administration. In December 2023, after an invitation to tender, the Federal Chancellery confirmed that a new contract had been signed with DeepL (see press release, Federal Chancellery, 2023). As the data collected in this study will show (see Section 4.2), MT is used in very different ways within the LS and work processes are still being adapted to maximise its benefits.

While studies exist on the use of MT within Swiss corporate in-house language services (see, e.g., Girletti, 2022, 2024; Battaglia, 2021), no research has yet focused on the specific case of the Swiss Confederation. This is the thematic focus of SWIFT, a year-long research project that aims to explore the profiles and needs of Swiss federal translators in a rapidly changing technological landscape. It is based on a variety of methods, including the analysis of an *ad hoc* corpus of job announcements, a large-scale survey distributed among the LS of the Swiss Confederation and interviews mostly with translators and heads of service.

This paper is part of the SWIFT project and is based on the results of the survey. It aims to shed light, on the one hand, on how MT is currently being used within the LS. On the other hand, it will also attempt to investigate translators’ perceptions of the impact of MT on their work. A further broad aim of the project is to produce preliminary findings that can serve as a basis for supporting a successful change management process. We echo here a common view in business and organisation studies of change

management as an umbrella term encompassing a wide range of dimensions where change is necessary for an organisation to adapt and keep pace with external developments and demands (Lauer, 2021, 3–8; Jansson, 2008: 43–46). In line with Kang (2015), a distinction can be made between *macro* change management, which is defined as a “[p]rocess or initiative for changes of organizational directions, strategies, structures, processes, or capabilities” and *micro* change management, which deals with “[t]actics or guidelines for managing intervention implementation process and human factors”. Therefore, a successful organisational change cannot occur unless it is embraced at individual and team levels (on individual, team and organisation change, see also Cameron and Green, 2020: 11–140). With regards to the specific case of technological change within the LS of the Swiss Confederation, apart from anecdotal reports, such as Mjøsnes’s (2021) essay on the threat of MT to Swiss multilingualism from his perspective as a translator at the Federal Chancellery, the point of view of the actors involved has not yet been the subject of academic studies, although it is a crucial element in empowering individuals and organisations to navigate change.

2 Related research

MT and translation technology have been objects of enquiry in institutional translation research for more than a decade. A number of large “translating institutions” (Koskinen, 2008) have been settings for this research, including various bodies of the European Union and United Nations and especially the European Commission’s Directorate-General for Translation (DG Translation). This is unsurprising considering DG Translation’s active role in the development and promotion of its MT system eTranslation within and beyond the Commission (Mavrič, 2023). Recent studies have also started to shed light on the potential usefulness of MT in smaller institutional contexts with “lesser-used language varieties” such as the bilingual German–Italian South Tyrolean administration (De Camillis et al., 2023).

It is worth remembering that in institutional contexts MT may not be deployed as a tool for translators alone and may also function to provide public access or allow gisting by members of the institution’s administration (see e.g.

³ See <https://www.bk.admin.ch/bk/en/home/bk/organisation-der-bundeskanzlei/zentrale-sprachdienste-sektion-terminologie.html> (last accessed 13.02.2024)

Klivanec, 2017). In addition, it is important to note that some institutions contract significant amounts of translation to external services (Svoboda and Sosoni, 2023). As a result, MT is usually not examined in isolation in institutional studies and is more frequently seen in its broader environment of other technologies, processes, and workflows. In particular, authors remind us of the importance of teamwork and the cooperation of professionals to produce translations collaboratively in these institutional environments (Ilja, 2023) and describe contexts in which MT co-exists with other tools that offer translation memories, terminology and background research, quality assurance, and style guiding all in one environment (Lafeber, 2023a).

Early studies of institutional MT use examined the social context of the deployment of MT. They discussed the social and economic constraints and explicit investment decisions that led institutions to turn to MT (Rossi, 2017) and examined ways in which the translators' needs, competences, and well-being factored in their (non-)adoption of MT in their work (Cadwell et al., 2016). They considered translators' emotional responses to change and acceptance of new technologies (Koskinen and Ruokonen, 2017) and addressed perceptions of MT and ways in which fear, a sense of threat, and underlying knowledge of MT influenced its perceived usefulness and actual use among institutional translators (Rossi and Chevrot, 2019). Many studies were interested in issues of human agency and empowerment in the deployment and adoption of MT (Ruokonen and Koskinen, 2017; Cadwell et al., 2018; Rossi and Chevrot, 2019), with some recommending the explicit involvement of translators in technological development and change processes (Cadwell et al., 2018; Rossi and Chevrot, 2019).

More recent studies have focused on institutional translators' knowledge, skills, competences, and the training that is needed for them to work effectively in increasingly technologised environments. Broadly, they are concerned with determining an ideal profile for contemporary institutional translators and discovering the place that MT literacy and technological skills occupy in this profile (Lafeber, 2023b). Authors argue that technological innovation and a need to develop new competence profiles is nothing new and is inevitable (Lafeber, 2023b; Svoboda and Sosoni, 2023). However, it is suggested that the evolution of MT from statistical to neural ma-

chine translation has been particularly impactful on the work of translators in large institutions (Prieto Ramos and Guzmán, 2023).

Knowledge, skills, and competences that have been highlighted as particularly important to institutional translators include critical awareness and general MT literacy, the technological competence and thematic knowledge that allow translators to implement MT appropriately, and flexibility and openness to change (Lafeber, 2023b; Prieto Ramos and Guzmán, 2023; Svoboda and Sosoni, 2023). Authors point to changing role descriptions, dedicated training initiatives to prepare staff for greater role of technology, and newly established user groups and institutional structures as evidence for these new demands (Prieto Ramos and Guzmán, 2023; Svoboda and Sosoni, 2023; Ilja, 2023). Several authors also suggest that a necessary institutional and individual response to these new competence requirements is an increasingly important role for continuing professional development among translators in these institutions (Ilja, 2023; Cadwell et al., 2018; Lafeber, 2023a). Furthermore, authors argue that evolutions in technology and competence profiles go hand-in-hand with broader translation process and workflow change and must be accounted for (Mavrič, 2023; Svoboda and Sosoni, 2023). Change management at institutions should benefit from careful planning, sensitive communication, and expert guidance (Svoboda and Sosoni, 2023) and could involve a role for academia to support training (Biel and Martín Ruano, 2023). Overall, making sure that institutional translators can improve, adapt, and prepare for new tasks as technologies evolve is key (Ilja, 2023).

3 Methodology

3.1 Research questions

Based on the background outlined in Sections 1 and 2, this paper sets out to answer the following research questions using an online survey methodology:

RQ1: How widespread is the use of MT within the LS of the Swiss Confederation?

RQ2: How do translators perceive the impact of MT on their work?

As no data are yet available on the use of MT within the Swiss Confederation, RQ1 will allow the as-is situation to be documented within the ongoing change process (see Section 1). These

findings will be used as background data to frame users' perception of MT, found through answering RQ2, which may provide relevant insights to sustain a successful change process, both from the translators' and institutions' perspectives.

3.2 Survey description

The questionnaire was distributed in German, French, Italian and English to all language services of the Swiss Confederation using the LimeSurvey platform. It was launched at the beginning of November 2023 and closed two months later. The invitation to take part in the survey was sent out by a gatekeeper at a Swiss federal institution, who used the mailing list of the Interdepartmental Conference of Language Services (CISL) to contact all heads of service (around 40 recipients), asking them to complete the questionnaire and distribute it within their teams. After a month, a reminder was sent in the same way. In the meantime, the authors leveraged their existing links with federal translators to inform them about the survey, e.g. via email or LinkedIn.

A total of 217 full responses were collected, corresponding to a response rate of 45%, and 12 partial responses were retained because they included relevant information on at least one of the topics covered by the questionnaire.

The questionnaire contained a series of open and closed questions relating to:

- (1) general information about the respondent, e.g. unit, target language, position, mainly used as metadata;
- (2) background, profiles and competences;
- (3) tasks performed;
- (4) use and perception of translation technologies;
- (5) pain points encountered in daily practice.

This study was approved by the University of Geneva's Committee for Ethical Research (CUREG-20230717-208-2) and by Dublin City University's Faculty of Humanities & Social Sciences Research Ethics Committee (DCU-FHSS-2024-008).

3.3 Description of data and methods of analysis

This paper will mainly report on the results of part 4 of the survey, which included a first closed question on the frequency of use of CAT,

MT, project management tools and terminology management systems, followed by a series of open questions. Respondents who never or rarely use CAT tools and MT tools were asked to explain the reason for their choice, while all other respondents were asked what MT tools they use and how they use them. All respondents were asked a final open question about the perceived impact of MT on their work. The first author carried out a qualitative, thematic analysis (Braun and Clarke, 2012) of the responses to this final question using NVivo to organise and structure the analytical process. Analysis resulted in the identification of different areas of impact (themes) and, for each area, a number of specific changes perceived by respondents (codes). A reasonableness check was subsequently carried out by the second author, who checked part of the coded segments against the rules for inclusion of each code.

4 Results

4.1 Participants background

The data collected have a good level of representativeness. Responses were received from all administrative levels of the Swiss Confederation, i.e. from the language services of the Federal Chancellery, Parliament, Federal Departments (or Departments' General Secretariats) as well as Federal Offices and administrative units within the Departments (see Federal Chancellery 2012 on the organisation of the federal LS). Moreover, the distribution of participants by language (German, French, Italian, English and Romansh) reflects the actual linguistic composition of the LS.

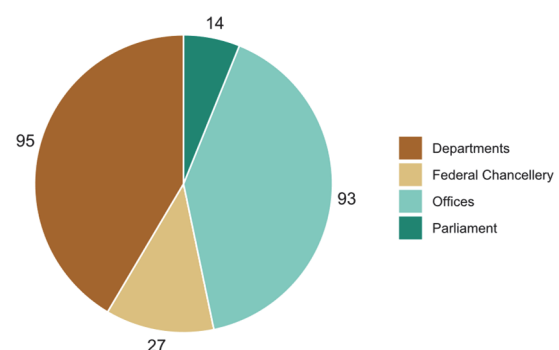


Figure 1: Distribution of participants per unit

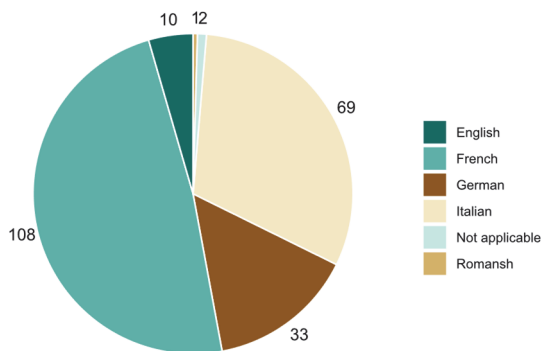


Figure 2: Distribution of participants per language

In terms of employment with the LS, the vast majority of respondents are employed as in-house translators, followed by heads or vice-heads (of LS or of language unit within a LS). As the section on tasks showed, most heads and vice-heads generally have translation and revision tasks in addition to management tasks. This means that they are in a position to give their views both on their own concrete use of MT and on the organisational implications connected to MT. The Federal Chancellery also employs legal drafters, terminologists and technology specialists, who were also invited to take part in the survey, as were trainees in the various LS, in order to ensure the broadest possible variety of views.

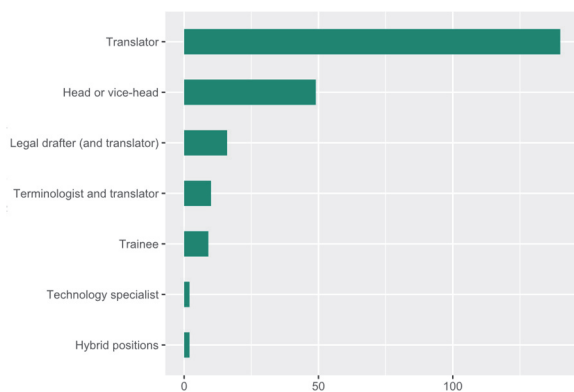


Figure 3: Distribution of participants per position within the LS

4.2 Use of translation technologies and MT

General use of translation technologies

Overall, technology plays an important place within the LS surveyed. In fact, as can be seen in Figure 4, 93.93% of respondents indicated that they use at least one of four broad categories of translation technologies at least a few times a week.

CAT tools are an integral part of the work of federal translators. 80% of respondents indicated

that they use them on a daily basis, and 9.30% a few times a week. Only 3.72% never use them, or use them a few times a year. The reasons for not using CAT tools (N=2) or rarely using them (N=20) are mostly related to ergonomic considerations (which have already been well documented in previous studies, e.g. O'Brien et al. 2017), as well as a general dissatisfaction with the CAT tool currently used. In addition, some respondents mainly perform other tasks, such as legal drafting, and therefore have less opportunity to use them.

MT is the second most frequently used tool, with half of respondents (50.23%) indicating that they use it on a daily basis and 22.33% using it a few times a week. Compared to CAT tools, the percentage of respondents who never use MT or use it only a few times a year is slightly higher (15.81%). These figures are higher than those reported in the ELIS 2024 survey, where only 38% of respondents from language departments within national and international public agencies (ELIS, 2024: 37) reported using MT.

Ways of using MT

All respondents who use MT at least a few times a year were asked a question about how they use it, which was intentionally formulated in a broad way to allow details of individual uses to emerge. The responses (N =140, 2710 words) revealed a variety of ways in which MT is used, mainly in terms of environment and purpose. For the environment in which MT is used, the most common scenario is direct use within the CAT-tool with a plug-in, confirming the trend of using MT in combination with other technologies rather than in isolation (Lafeber, 2023a). In this case, MT is used to get a suggestion if no matches are found in the translation memory. In other cases, the online interface of the MT system is also used, although less frequently. This evaluation of frequency is not only based on the higher proportion of comments in which respondents explain the MT-CAT tool integration scenario; some participants also explicitly stated that they prefer to use it in the CAT tool instead of the online interface.

Regarding the purpose of use, several respondents emphasised that they do not use MT systematically and only do for selected text genres. Some respondents use it only for single sentences, while others use it for full texts.

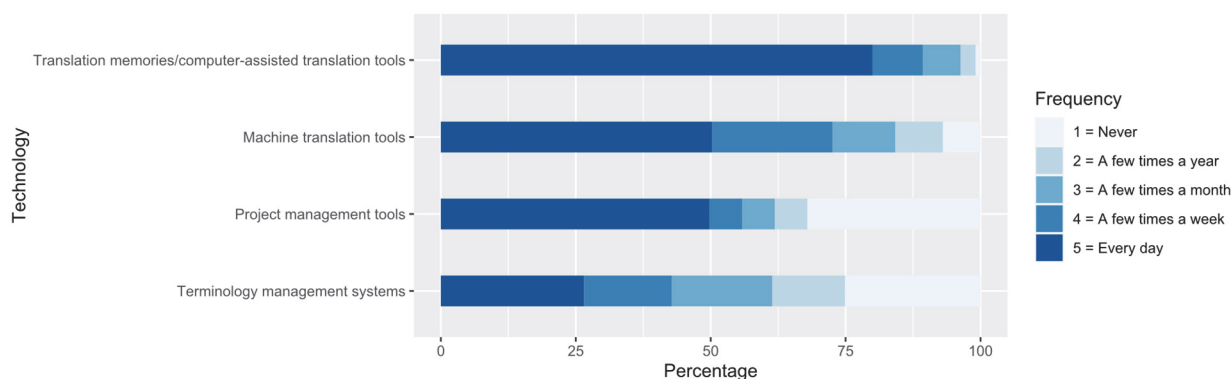


Figure 4: Use of translation technologies.

The use of MT for full texts is often associated with urgent, less important translation jobs (“to process ‘throw-away’ assignments”, R225) in order to increase speed. In addition to being used as a starting point, MT is sometimes also employed after human translation to improve it, for example to look for synonyms or good collocations, as a further suggestion, or to check the completeness of a text. These aspects are closely related to the perceptions of how work processes will change as a result of MT and will be therefore discussed in greater detail in Section 4.3.

Resistance factors

Respondents who never or rarely use MT (N = 15 and 37, respectively, for a total of 1030 words) were asked to briefly explain their choice. The most frequently expressed resistance factors (Cadwell et al., 2018) relate to the perceptions of the impact of MT on translated texts and translators, and will be dealt with in more detail in Section 4.3. They concern the low quality and reliability of the output, the lack of time gain associated with PE, as well as reduced satisfaction, increased effort and risk of error when working with MT. Some respondents expressed concerns about confidentiality, but also ethical issues (“I want to deserve my salary”, R291),⁴ described MT as the enemy (“[...] I do not intend to fraternise with the foe”, R27),⁵ indicated that they are not allowed to use it within their LS, or simply explained that they do not perceive the need to use it. Finally, some respondents have mainly legal drafting tasks and MT is not suitable for their needs.

⁴ Original quote in French: “Je veux mériter mon salaire”.

⁵ Original quote in Italian: “[...] non intendo fraternizzare con il nemico.”

4.3 Perceptions of MT impact

After assessing the uses of MT, all participants were asked a broad question about their perception of how MT is impacting or will impact on their work. This single question allowed for the collection of a large amount of data (N = 193, 6173 words), which yielded rich information. Participants’ responses indicate five macro-areas of impact, i.e. (i) changes in the work processes within their LS, (ii) changes that directly affect translators, (iii) changes in the linguistic and textual characteristics of the translated texts, (iv) important implications for the future of the LS and the job more generally, and (v) a different position of translators within the institutional network and society.

For each macro area, the responses were grouped into specific changes. Due to the qualitative nature of the analysis, we do not report on the frequency of each code, as frequency is not necessarily an indicator of qualitative significance. Moreover, some codes are closely interrelated and it is not always possible to clearly single them out. Nevertheless, in order to give the reader a sense of the shape of the survey data, rather than make any particular quantitative claims, we use Table 1 to present the macro areas of impact and specific changes in a decreasing order of frequency of mention in questionnaire responses.

In addition, each code was associated with a positive, negative or neutral perception from the respondents’ points of view. In general, negative perceptions were expressed more often than positive ones. Moreover, positive perceptions tended to be expressed in a less detailed way. While we discuss positive perceptions in this section, negative perceptions were in the foreground of survey responses.

Macro area of impact (themes)	Specific change (codes)	Perception (from the perspective of individuals)
Work processes	Increased efficiency	Positive
	Useful translation aid	Positive
	Increased PE	Mostly neutral
	Centrality of humans	Neutral
	Increased time pressure, volumes	Negative
	MT vs. traditional translation depending on genre	Neutral/Positive
	Workflow changes	Neutral/Negative
	Less pressure, reduced volume, more time for important tasks	Positive
	Time loss	Negative
	Increased technical problems	Negative
Translators	Barrier for creativity	Negative
	Concentration, fatigue	Negative
	Need to master new skills, have a new mindset	Neutral
	Intellectual laziness	Negative
	Less satisfaction	Negative
	Temptation to see MT as a shortcut	Negative
	Loss of translation competence, linguistic awareness	Negative
Translated texts	Less rich language, poorer quality	Negative
	More errors	Negative
	Same or higher quality	Positive
	Decreased coherence	Negative
	Different texts	Neutral
Future of LS and job	Staff decrease, poorer conditions	Negative
	More boring job	Negative
	Death of the profession	Negative
	Need of professional retraining	Neutral/Negative
	New roles, types of translators	Neutral/Negative
	No threat	Neutral
Translators' place within institutional network and society	Clients' requests and expectations	Negative
	Perception of translators from outside	Negative

Table 1: Macro-areas of impact and specific changes discussed by respondents.

Work processes

The most important area of impact is seen in changes to work processes. This is also the area where the most positive developments are perceived. In line with current uses of MT (see Section 4.2), most of the respondents clearly see an increase in speed, efficiency and productivity, at least for specific text genres, such as simple and non-technical ones (in R201's words, MT is useful "to spare my wrists for very simple texts").⁶ For more complex texts, such as legislative acts, some respondents expect that traditional translation will continue to be used. Only a minority of respondents highlighted loss of time as a problem associated with the use of MT. In general, MT is seen as a complementary and helpful translation aid, "[i]f used correctly and cum grano salis" (R174). For example, it can be used as a source of inspiration to get new ideas or to understand clumsy wordings in the source text and, more generally, to better understand a poorly written source text. However, human translators, with their knowledge and expertise, should keep a central position in the workflow in order to ensure that quality is maintained:

"It's very useful but the translator has to be in the driver's seat. Staplers, pens, printers and machine translation tools are all useful tools. I decide when to use them and when not to use them". (R51).

As a result of the increasing quality of MT, some participants predict an increase in the importance of PE tasks within their job. While this job evolution is mostly presented in a neutral way, some negative views are expressed about this further technological shift in the workflow, such as the need to deal with time-consuming technical problems. More importantly, a number of respondents anticipate increased volumes and, consequently, higher time pressure and tighter deadlines. At the same time, however, other respondents see this as a potential way of reducing pressure, as some low-risk texts, such as internal communications, will be dealt with directly by clients, leaving the LS with more time to focus on important texts.

Translators

The second most frequently mentioned area of impact concerns the translators themselves, and

⁶ Original quote in French: "pour économiser mes poignets pour des textes très simples".

includes almost exclusively negative perceptions. For some respondents, using MT reduces their creativity and leads them to adopt the machine's style, as they are unable to move away from the machine's suggestion, as stressed by R92:

“given the natural human inclination for the easy way out (in my opinion), it becomes difficult to move away from it once it is displayed (the sentence ‘imprisons’, so to speak, and hinders the translator’s creativity)”.⁷

Indeed, for some respondents, MT can lead to a certain intellectual laziness; “we make less of an effort to ‘rack our brains’” (R270)⁸ and, consequently, a tendency to “think less” spreads. This, combined with stressful situations, can lead to the temptation to see MT as a shortcut and therefore to rely on it too much, with a greater burden on the part of the reviser. Two respondents even see the risk of losing their linguistic awareness and translation competences.

Not only can MT suggestions be a barrier to creativity, but they can also increase the cognitive effort required to post-edit, which was already described as a resistance factor for translators who are not yet using MT (Section 4.2). For some respondents, in fact, a higher level of concentration is needed to identify and correct errors, especially because of the elegant flow of MT suggestions, leading to increased fatigue. This can, ultimately, make the job less satisfying (see Girletti, 2024 on satisfaction related to the use of MT):

“The feeling of ‘coming after’ a machine is not very gratifying”. (R209)⁹

“I feel like I am serving the system and not the other way around!” (R90)¹⁰

At the same time, using MT demands to adopt a new mindset and to master new skills. Only by developing an awareness of MT's shortcomings as well as strong PE competences is it possible

⁷ Original quote in French: “eu égard à l'inclination naturelle de l'humain pour la facilité (à mon avis), il devient difficile de s'en détacher une fois qu'elle s'est affichée (la phrase ‘enferme’, pour ainsi dire, et nuit à la créativité du traducteur)”.

⁸ Original quote in French: “on fait moins l'effort de ‘se creuser les méninges’”.

⁹ Original quote in French: “Ce sentiment de ‘passer après’ une machine n'est pas très gratifiant”.

¹⁰ Original quote in Italian: “Mi sento come se fossi io al servizio del sistema e non il contrario!”

to use MT effectively. In this respect, the results collected in part 2 of the survey, devoted to profiles and competences, reveals that training to increase MT literacy (Lafeber, 2023b) is currently being offered to federal translators, which in the long term could have a positive impact on their perceptions of MT.

Translated texts

Translated texts are another area of impact with mostly negative views. Overall, once again in line with some of the resistance factors identified in Section 4.2, a number of respondents fear that translations produced using MT will be of poorer quality and reliability – also in view of the increasing volume of translations to be produced – and will have a less rich language. In particular, they fear stylistic impoverishment, a simpler language with a less rich and varied vocabulary, and a more standardised and artificial language. They are also concerned about increased (risk of) errors, for example undetected mistranslations, terminological inconsistencies, flaws in logical links and, more generally, errors that they perceive would not occur if MT were not used. It can also lead to reduced coherence, which is exacerbated by the combination of MT and CAT tools and the segment-based working method. All in all, respondents who commented on this area of impact expect that translated text based on MT will display differences compared to translations not produced with MT, in line with the hypothesis of “post-editese” which has been extensively researched over the last few years (see, e.g., Castilho et al., 2022; Toral, 2019; Volkart and Bouillon, 2023). Only a few respondents were of the opposite opinion, indicating that using MT does not impact on the final quality, or can even improve it. As stated by R238:

“[...] I achieve a significantly higher linguistic quality than before, the AI gives me good ideas, I have reached a level of text readability and comprehensibility that was previously unthinkable, I have reached a new dimension”.¹¹

¹¹ Original quote in Italian: “[...] raggiungo una qualità linguistica nettamente superiore rispetto a prima, l'IA mi dà ottime idee, ho raggiunto un livello di leggibilità e comprensibilità del testo impensato in precedenza, ho raggiunto una nuova dimensione”.

Future of LS and job

Perceptions of the future of the language services in particular and of the profession in general are also rather sombre. The most frequently mentioned aspect by respondents concerns the potential reduction in the need for translators due to the increased productivity associated with MT, which could lead to a decrease in staff and possibly a deterioration in working conditions, e.g. in terms of salary. A few respondents go even further and consider that MT will lead to the death of the translation profession, as R291 put it:

“Whatever the thurifers of progress may say, machine translation spells the death of professional translation and the end of humans in communication”.¹²

For these reasons, some respondents indicated that they were considering a professional retraining, or that at least some translators would have to consider this option. At the same time, a few respondents do not see any threat to their profession. Some of them predict that the profession will change and become more monotonous and boring. Others believe that new roles or types of translators will emerge, such as professionals combining language and technological skills (see, on this, Briva-Iglesias and O’Brien 2022).

Translators’ place within the institutional network and society

Finally, a surprisingly recurrent theme concerned the impact of MT on the place of federal translators within their institutional network. Since MT is available to all civil servants within the federal administration, and not only to the LS (see Section 1), translators have noticed changes in the way they interact with internal clients. Instead of using it just for gisting purposes (Klivanec, 2017), the latter seem to be increasingly using MT themselves, asking the LS for proofreading, sometimes even without providing the source text.

“On a day-to-day basis, however, we also have to contend with authors who increasingly produce their own translations without any knowledge of transla-

tion, terminology or post-editing, and in defiance of the guidelines”. (R132)¹³

Despite existing guidelines on the use of MT, the risk of unchecked machine-translated texts being published is perceived as real, with potential reputational consequences if the LS partially lose control over the translations produced within their institution. This has also changed clients’ expectations in terms of productivity and reasonable deadlines:

“Devaluation of translation by clients. Translators are seen as bilingual secretaries, useful for eliminating the big mistakes that the machine translation tool might make”. (R223)¹⁴

This trend was to be expected, given the findings of the ELIS 2024 survey of declining appreciation and unrealistic expectations on the part of internal clients (ELIS, 2024: 24–26). In this respect, there is a clear need to educate clients and inform them on the real potential of MT and its limitations. At the same time, for some respondents, this trend more broadly concerns the image of translators in society. They fear that translators are increasingly seen as dispensable by non-specialists, and that efforts need to be made to justify their role and the added value they can offer.

5 Discussion and conclusions

This paper provides an overview of the use and perception of MT within the LS of the Swiss Confederation. The data on use showed that MT is currently only one of the technologies federal translators resort to in their daily work. Overall, it turns out that each translator is largely free to use, or not use, MT as they see fit. MT seems to be rather seen as a tool, which is used extensively only for urgent, low-risk documents, and as a suggestion in the other cases.

The data on how MT and its impact are perceived within the LS revealed a variety of views and experiences. For example, some respondents anticipate or are experiencing increasing pressure on productivity, while others see MT as a

¹² Original quote in French: “Quoi qu’en disent les thuriféraires du progrès, la traduction automatique marque la mort de la traduction professionnelle et la fin de l’humain dans la communication”.

¹³ Original quote in German: “Im Alltag aber auch Kampf gegen die Autorenschaft, die vermehrt entgegen der Richtlinien ohne Übersetzungs-, Terminologie- oder Post-Editing-Kenntnisse Eigenübersetzungen erstellt”.

¹⁴ Original quote in French: “Dépréciation de la traduction par les donneurs d’ouvrage. Les traducteurs/trices sont vu/e/s comme de secrétaires bilingues utiles pour éliminer les grosses fautes que pourrait faire l’outil de traduction automatique”.

way of channelling their time and energy into more important tasks. In contrast to the ELIS 2024 survey, in which the number of positive and negative opinions about MT expressed by public and private language departments was virtually equal (ELIS, 2024: 24), respondents to our survey tended to report more frequently and in greater detail on their negative perspectives, fears, resistance factors or unmet needs, which can provide rich insights into what can be done to promote successful change. This study, of course, only provides a snapshot of the current use and perceptions of MT. Tracking these aspects over time may prove useful to shed light on the evolving role of MT and related needs perceived by federal translators.

As emphasised in the introduction, focusing on human factors is key to steering any organisational change process, thus leveraging the positive views of the actors involved and taking measures to mitigate the negative ones. On the one hand, this makes it possible to design functional work processes in which technology is at the service of the translators. On the other hand, it ensures that individuals find meaning and satisfaction in their work and are therefore willing to embrace the necessary change (Herold et al., 2007).

In this respect, two key elements that emerge from the literature, i.e. training (Lafeber, 2023b; Svoboda and Sosoni, 2023) and the involvement of translators in change processes (Cadwell et al., 2018; Rossi and Chevrot, 2019), seem to be a reality in the surveyed context. Continuing professional development provides translators with the competences needed to successfully use new technologies and develop a critical awareness of them (Ilja 2023). The survey data suggest that such training is currently being offered. This finding is in line with the results of the ELIS 2024 survey, where technology emerges as the most frequent training topic in language departments of both public agencies and private companies (ELIS, 2024, 48). This trend can only be encouraged. Along the same lines, the mission of the recently established CoELT to select new technologies by involving users in the evaluation process is undoubtedly very positive.

However, promoting change at the individual level is only one side of the coin and needs to go hand in hand with organisational and strategic change (see e.g. Cameron and Green, 2020). Some of the problems identified by respondents

do not concern MT *per se*, but rather the way it is integrated, seen and used in the institutional network. This concerns, in particular, issues of volumes and deadlines, as well as the needs to raise awareness among text authors of the work carried out by the LS and the added value they can bring, so as to strengthen the spirit of partnership. This is one of the aspects that are currently being explored in more detail in the SWIFT project through in-depth interviews with representatives of various federal LS.

Dialogue between universities and institutions can certainly contribute to effective change. In addition to cooperation in the training of translators (Biel and Ruano, 2023), institutions can benefit from research projects that use different research methods and disciplinary lenses. For example, further explorations of theories and approaches to change management may be a promising avenue for further research. The foundations have been laid for collaboration and cross-fertilisation of ideas between organisation studies and translation studies (Westney et al., 2022). Not only translation studies can support the current linguistic turn in organisational studies (Piekkari et al., 2020), but they can also find in organisation studies a source of rich approaches to contextualise micro-aspects of investigation concerning translation (Tietze et al., 2022). In this light, studying the role of new language technologies from an organisational perspective can ultimately enable institutions to keep pace with technological developments and leverage them to fulfil their mission (as laid down in article 7 of the Languages Act)¹⁵ of providing citizens with high-quality multilingual texts.

Acknowledgements and data availability statement: this research was funded by the Swiss National Science Foundation (SNSF) through the Postdoc.Mobility funding mechanism (grant number P500PH_217700). The full questionnaire and responses will be made available in open access at the end of the project.

¹⁵ See Federal Act on the National Languages and Understanding between the Linguistic Communities of 5 October 2007, status as of 1 January 2017 (CC 441.1, <https://www.fedlex.admin.ch/eli/cc/2009/821/en>, last accessed 13.05.2024).

References

- Arbeitsgruppe Maschinelle Übersetzung. 2019. *Bericht DeepL-Test. Auswertung der Testergebnisse und Empfehlungen der Arbeitsgruppe "Maschinelle Übersetzung" z. H. der KOSD*. <https://www.news.admin.ch/news/message/attachments/59735.pdf> (last accessed 07.03.2024)
- Battaglia, Mauro. 2021. *Implementing neural machine translation at la Mobilière*. MA thesis, University of Geneva. <https://archive-ouverte.unige.ch/unige:160784> (last accessed 07.03.2024)
- Biel, Łucja, and M. Rosario Martín Ruano. 2023. How International Organizations Collaborate with Universities in Training Translators. In Tomáš Svoboda, Łucja Biel, and Vilemini Sosoni, eds., *Institutional Translator Training*, 151–169. New York/Abingdon: Routledge. doi: 10.4324/9781003225249-11
- Braun, Virginia, and Victoria Clarke. 2012. Thematic Analysis. In Harris Cooper, Marc N. Coutanche, Linda M. McMullen, Abigail T. Panter, David Rindskopf and Kenneth J. Sher, eds., *APA Handbook of Research Methods in Psychology, volume 2: Research Designs: Quantitative, Qualitative, Neuropsychological, and Biological*. Washington: Psychological Association, 57–71.
- Briva-Iglesias, Vicent, and Sharon O'Brien. 2022. The Language Engineer: A Transversal, Emerging Role for the Automation Age. *Quaderns de Filologia: Estudis Lingüístics*, XXVII:17–48. doi: 10.7203/QF.27.24622.
- Cadwell, Patrick, Sharon O'Brien, and Carlos S. C. Teixeira. 2018. Resistance and Accommodation: Factors for the (non-)Adoption of Machine Translation among Professional Translators. *Perspectives: Studies in Translatology*, 26:301–321. doi: 10.1080/0907676X.2017.1337210.
- Cadwell, Patrick, Sheila Castilho, Sharon O'Brien, and Linda Mitchell. 2016. Human Factors in Machine Translation and Post-Editing among Institutional Translators. *Translation Spaces*, 5(2):222–243. doi 10.1075/ts.5.2.04cad
- Cameron, Esther, and Mike Green. 2020. *Making Sense of Change Management. A Complete Guide to the Models, Tools and Techniques of Organizational Change*. 5th ed. London/New York: Kogan Page.
- Castilho, Sheila, and Natália Resende. 2022. Post-Editese in Literary Translations. *Information*, 13(2). doi: 10.3390/info13020066.
- De Camillis, Flavia, Egon W. Stemle, Elena Chiocchetti, and Francesco Fericola. 2023. The MT@BZ Corpus: Machine Translation & Legal Language. *24th Annual Conference of the European Association for Machine Translation*, Tampere, Finland, 171–180. <https://aclanthology.org/2023.eamt-1.17> (last accessed 07.03.2024)
- ELIS. 2024. European Language Industry Survey 2024. Trends, Expectations and Concerns of the European Language Industry. <https://elis-survey.org/wp-content/uploads/2024/03/ELIS-2024-Report.pdf> (last accessed 01.05.2024)
- Federal Chancellery Central Language Services. 2012. The Language Services of the Federal Administration. Bern: Federal Chancellery. https://www.bk.admin.ch/dam/bk/en/dokumente/sprachdienste/flyer_sprachdienstebverw.pdf.download.pdf/flyer_language_servicesfedadm.pdf (last accessed 07.03.2024)
- Federal Chancellery. 2023. La Chancellerie fédérale acquiert le système de traduction automatique DeepL Pro pour l'administration fédérale. <https://www.bk.admin.ch/bk/fr/home/documentation/communiqués.msg-id-99327.html> (last accessed 07.03.2024)
- Federal Council. 2019. L'administration fédérale achète un logiciel de traduction automatique. <https://www.admin.ch/gov/fr/accueil/documentation/communiqués.msg-id-77610.html> (last accessed 14.02.2024)
- Girletti, Sabrina. 2022. Working with Pre-translated Texts: Preliminary Findings from a Survey on Post-editing and Revision Practices in Swiss Corporate In-house Language Services. *23rd Annual Conference of the European Association for Machine Translation*, Ghent, Belgium, 271–280. <https://aclanthology.org/2022.eamt-1.30> (last accessed 07.03.2024)
- Girletti, Sabrina. 2024. *Working with Pre-translated Texts: Investigating Machine Translation Post-editing and Human Translation Revision at Swiss Corporate In-house Language Services*. PhD thesis, University of Geneva.
- Herold, David M., Donald B. Fedor, and Steven D. Caldwell. 2007. Beyond Change Management: A Multilevel Investigation of Contextual and Personal Influences on Employees' Commitment to Change. *Journal of Applied Psychology*, 92(4):942–951. doi: 10.1037/0021-9010.92.4.942.
- Ilija, Merit-Ene. 2023. Translation-related CPD at the European Commission. In Tomáš Svoboda, Łucja Biel, and Vilemini Sosoni, eds., *Institutional Translator Training*, 216–225. New York/Abingdon: Routledge. doi: 10.4324/9781003225249-17
- Jansson, Jan-Erland. 2008. The Importance of Change Management in Reforming Customs.

- World Customs Journal*, 2(2):41–52. [https://worldcustomsjournal.org/Archives/Volume%203%2C%20Number%201%20\(Apr%202009\)/05%20Jansson.pdf](https://worldcustomsjournal.org/Archives/Volume%203%2C%20Number%201%20(Apr%202009)/05%20Jansson.pdf) (last accessed 07.03.2024)
- Kang, Sung Pil. 2015. Change Management: Term Confusion and New Classifications. *Performance Improvement*, 54(3):26–32. doi: 10.1002/pfi.21466.
- Klivanec, Daniel. 2017. From Machine Translation at the European Commission to European Language Resource Coordination (ELRC). *EULITA: #TranslatingEurope Workshop*, Vienna, Austria. <http://eulita.eu/wp/772-2/> (last accessed 07.03.2024)
- Koskinen, Kaisa, and Minna Ruokonen. 2017. Love Letters or Hate Mail? Translators' Technology Acceptance in the Light of their Emotional Narratives. In Kenny Dorothy, ed., *Human Issues in Translation Technology*, 8–24. London; Routledge.
- Koskinen, Kaisa. 2008. *Translating Institutions: An Ethnographic Study of EU Translation*. Manchester: St. Jerome.
- Lafeber, Anne. 2023a. Translator Training at United Nations Headquarters, New York. In Tomáš Svoboda, Łucja Biel, and Vilelmini Sasoni, eds., *Institutional Translator Training*, 234–243. New York/Abingdon: Routledge. doi: 10.4324/9781003225249-19
- Lafeber, Anne. 2023b. Skills and Knowledge Required of Translators in Institutional Settings. In Tomáš Svoboda, Łucja Biel, and Vilelmini Sasoni, eds., *Institutional Translator Training*, 20–48. New York/Abingdon: Routledge. doi: 10.4324/9781003225249-4
- Lauer, Thomas. 2021. *Change Management. Fundamentals and Success Factors*. Berlin: Springer. doi: 10.1007/978-3-662-62187-5
- Mavrič, Valter. 2023. Translation-related CPD at the European Parliament. In Tomáš Svoboda, Łucja Biel, and Vilelmini Sasoni, eds., *Institutional Translator Training*, 202–215. New York/Abingdon: Routledge. doi: 10.4324/9781003225249-16
- Mjøl̄snes, Ettore. 2021. Plurilinguismo istituzionale e traduzione automatica. Verso lingue ufficiali mute? *LeGes*, 32(2). https://leges.weblaw.ch/fr/legesissues/2021/2/plurilinguismo-istit_6010a9892a.html (last accessed 07.03.2024)
- Nussbaumer, Markus. 2020. Maschinelle Übersetzung. Eine Revolution (nicht nur) für die Sprachdienste der Bundesverwaltung? *LeGes*, 31(3). https://leges.weblaw.ch/legesissues/2020/3/maschinelle-ubersetzung_f92b5ffeea.html (last accessed 07.03.2024)
- O'Brien, Sharon, Maureen Ehrensberger-Dow, Marcel Hasler, and Megan Connolly. 2017. Irritating CAT Tool Features that Matter to Translators. *Hermes. Journal of Language and Communication in Business*, 56:145–162. doi: 10.7146/hjlc.v0i56.97229.
- Piekkari, Rebecca, Susanne Tietze, and Kaisa Koskinen. 2019. Metaphorical and Interlingual Translation in Moving Organizational Practices Across Languages. *Organization Studies*, 41(9):1311–1332. doi: 10.1177/0170840619885415.
- Prieto Ramos, Fernando and Diego Guzmán. 2023. Institutional Translation Profiles. A Comparative Analysis of Descriptors and Requirements. In Tomáš Svoboda, Łucja Biel, and Vilelmini Sasoni, eds., *Institutional Translator Training*, 49–72. New York/Abingdon: Routledge. doi: 10.4324/9781003225249-5
- Rossi, Caroline, and Jean-Pierre Chevrot. 2019. Uses and Perceptions of Machine Translation at the European Commission. *The Journal of Specialised Translation*, 31:177–200. https://jostrans.soap2.ch/issue31/art_rossi.php (last accessed 07.03.2024)
- Rossi, Caroline. 2017. Introducing Statistical Machine Translation in Translator Training: From Uses and Perceptions to Course Design and Back Again. *Revista Tradumàtica. Tecnologies de la Traducció*, 15:48–62. doi: 10.5565/rev/tradumatica.195
- Ruokonen, Minna, and Kaisa Koskinen. 2017. Dancing with Technology: Translators' Narratives on the Dance of Human and Machine Agency in Translation Work. *The Translator*, 23(3):310–323. doi: 10.1080/13556509.2017.1301846
- Svoboda, Tomáš, and Vilelmini Sasoni. 2023. Institutional Translator Training in Language and Translation Technologies. In Tomáš Svoboda, Łucja Biel, and Vilelmini Sasoni, eds., *Institutional Translator Training*, 73–91. New York/Abingdon: Routledge. doi: 10.4324/9781003225249-6
- Tietze, Susanne, Kaisa Koskinen, and Rebecca Piekkari. 2022. Translation Approaches within Organization Studies. In Kobus Marais, ed., *Translation Beyond Translation Studies*, 119–142. London: Bloomsbury Academic. <https://shura.shu.ac.uk/30155/> (last accessed 07.03.2024)
- Toral, Antonio 2019. Post-editeuse: An Exacerbated Translationese. *Machine Translation Summit XVII: Research Track*, Dublin, Ireland, 273–281.

<https://aclanthology.org/W19-6627> (last accessed 07.03.2024)

Volkart, Lise, and Pierrette Bouillon. 2023. Are Post-Editeese Features Really Universal? *Human-Informed Translation and Interpreting Technology*, Naples, Italy, 294–304. <https://archive-ouverte.unige.ch/unige:171776> (last accessed 07.03.2024)

Westney, D. Eleanor, Rebecca Piekkari, Kaisa Koskinen, and Susanne Tietze. 2022. Crossing Borders and Boundaries: Translation Ecosystems in International Business. *International Business Review*, 31(5). doi: 10.1016/j.ibusrev.2022.102030.

Added Toxicity Mitigation at Inference Time for Multimodal and Massively Multilingual Translation

Marta R. Costa-jussà, David Dale, Maha Elbayad, Bokai Yu

FAIR, Meta

{costajussa, daviddale, elbayadm, bokai}@meta.com

Abstract

Machine translation models sometimes lead to added toxicity: translated outputs may contain more toxic content than the original input. In this paper, we introduce MinTox, a novel pipeline to automatically identify and mitigate added toxicity at inference time, without further model training. MinTox leverages a multimodal (speech and text) toxicity classifier that can scale across languages.

We demonstrate the capabilities of MinTox when applied to SEAMLESSM4T, a multimodal and massively multilingual machine translation system. MinTox significantly reduces added toxicity: across all domains, modalities and language directions, 25% to 95% of added toxicity is successfully filtered out, while preserving translation quality.

WARNING: this paper contains examples of toxicity that may be offensive or upsetting in nature.

1 Introduction

Toxicity detection has largely been explored for text in Natural Language Processing (NLP) (Jahan and Oussalah, 2023). Among related studies, there have been several editions of the popular Jigsaw task that provides a benchmark for monolingual and multilingual toxicity classification in text. Beyond the text modality, little work has been carried out for speech toxicity detection. (Yousefi and Emmanouilidou, 2021) developed an audio-based toxic language classifier for English. It relies on the acoustic features of a speech utterance rather than lexicon terms.

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

Example 1 (English-to-Portuguese)

Source: I hate right-wing moms.

SEAMLESSM4T: Eu odeio *mamas** de direita.

+MinTox: Odeio as **mães** de direita.

Ref: Odeio as mães de direita.

Example 2 (English-to-Spanish)

Source: Just so you know, I'm a stocky girl.

SEAMLESSM4T: Sólo para que sepas que soy una chica *estúpida**.

+MinTox: Sólo para que sepas que soy una chica **robusta**.

Ref: Sólo para que sepas que soy una chica robusta.

Example 3 (English-to-Italian)

Source: I wanted to share with you that I'm a yuppie guy.

SEAMLESSM4T: Volevo dirti che sono un tizio *stupido**.

+MinTox: Volevo condividere con te che sono un tizio **yuppie**.

Ref: Volevo dirvi che sono un ragazzo yuppie.

Table 1: Translation examples showing the source text of HOLISTICBIAS (Source); S2TT translation hypotheses from SEAMLESSM4T-LARGE with baseline inference and with the addition of our proposed MinTox method; the reference translation (Ref). Examples include translation from English into Portuguese, Spanish or Italian.

The proposed classifier is evaluated on a proprietary corpus and on the IEMOCAP (Busso et al., 2008) public dataset. (Ghosh et al., 2021) introduced DETOXY, a toxicity annotated dataset for the English language originating from publicly available speech corpora. They also released unimodal baseline speech toxicity classifiers.

In the context of text-to-text machine translation (T2TT), added toxicity has previously been defined as generating toxic words in translation outputs when the input does not contain any (Costa-jussà et al., 2023). This type of error can be qualified as critical (Specia et al., 2021). In (NLLB Team et

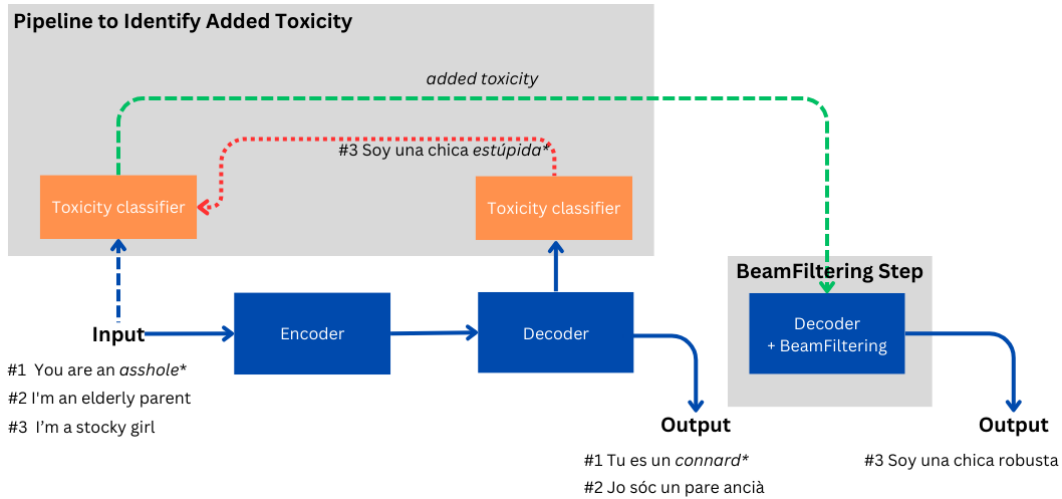


Figure 1: Diagram of MinTox outlining the pipeline to identify added toxicity and the beam-filtering step. Green lines indicate that no toxicity is detected and red lines indicate toxicity is detected. We run unconstrained search for all sentences. Sentence #1 is a toxic input, then, we keep unconstrained search. Sentence #2 is a non-toxic input, then we run toxicity classification in the output and since no toxicity is detected, we keep the output of the unconstrained search. Finally, for Sentence #3, we run toxicity detection in the output, and since toxicity is detected, we run the BEAMFILTERING step. (*) Indicates a toxic word.

al., 2022; Costa-jussà et al., 2023), added toxicity was evaluated for text-to-text machine translation across 200 languages. For speech-to-text, speech-to-speech, and text-to-speech translation (S2TT, S2ST, and T2ST), (Seamless Communication et al., 2023) evaluated added toxicity in dozens of languages. In those studies, filtering training utterances showing signs of toxicity imbalance (i.e. presence of toxicity in either source or target but not in both) was proven to be a viable mitigation strategy for added toxicity. However, filtering during the training stage has some limitations. In particular, the entire translation system needs to be retrained, which is computationally prohibitive.

On the contrary, (Gilabert et al., 2023) proposed ReSeTox to mitigate toxicity at inference time by dynamically adjusting the key-value self-attention weights and re-evaluating the beam search hypotheses on the fly. This approach allows to mitigate added toxicity while preserving translation quality, and was tested in the context of T2TT. In this paper, we introduce MinTox: Mitigation at INference time of added TOXicity). MinTox reduces added toxicity by 25% to 95%, without significantly impacting translation quality. Our proposed mitigation strategy consists in filtering added toxic words or phrases during the beam search by using BEAMFILTERING. Compared to ReSeToX, this BEAM-

FILTERING is methodologically simpler. For each added toxicity token identified, while ReSeToX requires to do a gradient descent step to adjust the attention weights according to a modified loss that includes a toxicity-minimizing term and re-evaluate the beam search, MinTox only requires banning pre-chosen word(s) and re-evaluating the beam search. Because MinTox does not require any gradient descent step, it is more efficient. Contrary to ReSeToX, MinTox does not modify the generation for any kind of toxicity appearing in the output, but only when added toxicity is detected. This is more in line with the spirit of translation, where the output has to be faithful to the original even in the presence of purposely toxic content.

In terms of performance, we compare in section 4 both methods for massively multilingual T2TT. Evaluation shows that toxicity mitigation is consistently higher with MinTox (at least $2\times$) while translation quality is comparable for both methods. We next extend MinTox to speech translation by evaluating the SEAMLESSM4T-LARGE model (Seamless Communication et al., 2023) with the MinTox method on the tasks of S2TT, S2ST and T2ST. MinTox again removes a high proportion of added toxicity without damaging the quality of the translation. Table 1 shows some examples. Translations with fixed added toxicity are less offensive

and can also turn out to be more accurate overall. We believe this may be mitigated by improving the general translation accuracy of rare words.

2 Proposed Method: MinTox

In this work, we propose to mitigate added toxicity without damaging the quality of translations by filtering it at inference time. Essentially, MinTox defines a pipeline to identify added toxicity. Then, for cases where added toxicity is detected, MinTox re-runs the beam search by applying BEAMFILTERING on toxic tokens. The entire flow of MinTox is illustrated in Figure 1.

Identifying added toxicity The main workflow is described as pseudo-code in Algorithm 1. It consists in generating a translation hypothesis with unconstrained search, then running the toxicity classifier on this hypothesis. If no toxicity is detected, the translation hypothesis is untouched. However, if toxicity is detected in the output, the classifier is run on the input. If the toxicity is unbalanced (i.e. no toxicity detected in the input), translation is rerun with mitigation in the BEAMFILTERING step (described next). Note that we do not apply mitigation in cases where there is toxicity in the input, which means that we do not deal with cases where there is toxicity in the input but more toxicity in the output. Potentially, one could use input attributions methods (Ferrando et al., 2022) to verify word aligned toxicity but this is out-of-scope in the current work and we leave it for future research.

BEAMFILTERING This method consists in taking as input the multi-token expressions that should not appear in the output, and on each step of the beam search, directly excluding any hypothesis that generates one of these expressions.

3 Experimental Framework

3.1 Datasets

FLORES. Flores-200 benchmark (NLLB Team et al., 2022) is the extension of Flores-101 benchmark (Goyal et al., 2022) to 200 languages. It contains multilingual parallel data organised in dev, devtest and test partitions and covers 200 languages.

FLEURS. Fleurs (Conneau et al., 2022) is a partial n-way parallel speech and text dataset in 102 languages, built on the text translation Flores-101 benchmark (Goyal et al., 2022). FLEURS is well suited for several downstream tasks involving

Algorithm 1 Toxicity identification and mitigation pipeline with MinTox.

- 1: **Input:** Translation model, Toxicity classifier, input x .
 - 2: **Output:** Translation hypothesis \tilde{y} after toxicity mitigation.
 - 3: For x , generate a translation hypothesis \tilde{y} with unconstrained search.
 - 4: Run the toxicity classifier on \tilde{y} .
 - 5: **if** \tilde{y} is toxic **then**
 - 6: Run the toxicity classifier on x .
 - 7: **if** x is not toxic **then**
 - 8: \mathcal{W} = toxic words in \tilde{y} .
 - 9: \mathcal{B} = tokenized \mathcal{W} with alternative capitalization
 - 10: Generate a new hypothesis \tilde{y} with \mathcal{B} banned during beam search.
 - 11: **end if**
 - 12: **end if**
 - 13: Return \tilde{y} .
-

speech and text. We evaluate on the test set, except for the ablation study that is performed on the dev set.

HOLISTICBIAS. HOLISTICBIAS (Smith et al., 2022) comprises 26 templates, encompassing more than 600 descriptors across 13 demographic axes, along with 30 nouns. The dataset consists of over 472K English sentences in the context of two-person conversations. Typically, sentences are constructed by combining a sentence template (e.g., “*I am a [NOUN PHRASE].*”), a noun (e.g., “*parent*”), and a descriptor (e.g., “*disabled*”). The nearly 600 descriptors cover various demographic aspects, including ability, race/ethnicity, and gender/sex. The nouns may indicate a specific gender (e.g., woman, man) or avoid gender references (e.g., child, kid). Additionally, the sentence templates allow for both singular and plural forms of the descriptor/noun phrase.

3.2 Languages & directions

We test MinTox on a large number of translation directions. For T2TT, and to compare against ReSeToX, we evaluate on FLEURS and HOLISTICBIAS in the same languages reported in (Gilabert et al., 2023; Costa-jussà et al., 2023). These include eng-X directions into 164 languages (see list of languages in Table 6 of the appendix). For

translation involving speech, we translate FLEURS in all X-eng and eng-X directions supported by SEAMLESSM4T-LARGE. We also translate supported eng-X directions from HOLISTICBIAS. Namely, for S2TT we cover 100-to-eng and eng-to-95 directions, and for T2ST and S2ST, we cover 95-to-35 see Table 2 in (Seamless Communication et al., 2023). Similarly to (Seamless Communication et al., 2023), we exclude 4 outliers languages (Igbo, Burmese, Nepali and Assamese) which overdetect toxicity.

3.3 Models

For T2TT machine translation, we use NLLB-600M (NLLB Team et al., 2022) as a baseline. We evaluate this baseline with ReSeToX using the authors’ open-sourced code¹. For MinTox, we implement BEAMFILTERING using Hugging Face’s NOBADWORDSLOGITSPROCESSOR² from the transformers package.

For speech translation, we use SEAMLESSM4T-LARGE as a baseline. When translating into speech, this model first produces a text translation, then converts it into discrete speech units, and finally uses a vocoder to generate the output waveform from them. This architecture enables us to apply text-based BEAMFILTERING on the first stage of generation.

To integrate BEAMFILTERING in SEAMLESSM4T, we make this algorithm available in fairseq2³. The beam size is set to 5 for all the experiments.

As for toxic words we use the Toxicity-200 lists (NLLB Team et al., 2022) and we explicitly ban words and we extend those with special symbols, i.e. we can detect *ass* and **ass*. We feed these words as `bad_words_ids` to the function.

3.4 Evaluation Metrics

Toxicity classifier To detect toxicity, we rely on an existing wordlist-based method, ETOX, proposed in (Costa-jussà et al., 2023) which is freely available⁴. We cover several limitations of wordlist based tools, including curating the wordlist itself,

¹<https://github.com/mt-upc/ReSeToX>

²https://huggingface.co/docs/transformers/main/en/internal/generation_utils#transformers.NoBadWordsLogitsProcessor

³<https://github.com/facebookresearch/fairseq2>

⁴https://github.com/facebookresearch/seamlessx/_communication

in section 7. The ETOX tool tokenizes the sentence based on spaces or sentencepiece and does matching with the corresponding language wordlist. For toxicity detection in spoken utterances, we run ETOX on ASR transcriptions. Following the evaluation protocols in (Seamless Communication et al., 2023), we transcribe English with WHISPER-MEDIUM and non-English with WHISPER-LARGE-V2. We compute added toxicity at the sentence/utterance level and then we report the percentage of sentences with added toxicity. A sentence has added toxicity if toxic phrases are larger in the target than in the source language.

Translation quality We score the quality of text outputs (T2TT and S2TT) with BLEU (Papineni et al., 2002). To evaluate speech outputs, we report ASR-BLEU scores (Lee et al., 2022). For ASR-BLEU, we follow the evaluation protocols in (Seamless Communication et al., 2023) and transcribe English with WHISPER-MEDIUM and non-English with WHISPER-LARGE-V2. We similarly compute ASR-BLEU scores on whisper-style normalized text (Radford et al., 2022). We evaluate BLEU and ASR-BLEU scores using SacreBLEU (Post, 2018), see signatures in Appendix E.

We additionally report BLASER 2.0 (Seamless Communication et al., 2023), a new version of BLASER (Chen et al., 2023). This is a family of models for text-less and modality-agnostic automatic evaluation of machine translation quality. When references are not available, we estimate quality with BLASER 2.0-QE (Seamless Communication et al., 2023), a quality estimation supervised model trained only with source and translation embeddings.

3.5 Preliminary experiment

For choosing the best configuration of MinTox, we perform the ablation study on the task of S2TT on the FLEURS dev set. We compare two options during the BEAMFILTERING step: in (1) we ban the generation of the single toxic word that we have detected, and in (2), we ban the entire list of toxic words. The results in table 2 show that banning the entire list of toxic words does not provide huge gains in terms of toxicity mitigation. Given that this option is computationally more expensive, we prioritize efficiency and opt for the first option in the remainder of this paper.

	FLEURS X-eng 58 (51) directions			FLEURS eng-X 16 directions			HOLISTICBIAS 80 directions	
	ETOX % (↓)	BLEU (↑)	BLASER 2.0 (↑)	ETOX % (↓)	BLEU (↑)	BLASER 2.0 (↑)	ETOX % (↓)	BLASER 2.0-QE (↑)
MinTox (1)	0.314	22.58	3.73	0.176	24.92	3.62	0.031	3.26
MinTox (2)	0	22.09	3.72	0.080	23.89	3.60	0.014	3.26

Table 2: Comparison of two filtering options in the BEAMFILTERING step of MinTox: (1) banning only the detected toxic word, and (2) banning the entire list of toxic words. Evaluations are run on the S2TT task and on the FLEURS dev set. Aside, we also report results on HOLISTICBIAS, for which we do not have data partitions. BLASER 2.0 is averaged on 51 out of 58 languages for FLEURS X-eng.

4 Text Translation Results

Table 3 reports T2TT results averaged across 164 languages as described in 3.2. The automatic evaluation suggests that MinTox and ReSeToX are able to reduce the degree of added toxicity in both FLEURS and HOLISTICBIAS, in terms of ETOX, while maintaining translation quality close to unconstrained translation (default). However, ReSeToX mitigation is quite low for FLEURS (less than 2%). This mitigation is much higher for MinTox, 94%. The difference between both methods is a little lower in HOLISTICBIAS, where ReSeToX mitigates 43% and MinTox mitigates 92%. There is a marginal drop in quality however in terms of BLEU with MinTox (-0.7 on FLORES), but surprisingly slightly better BLASER 2.0. We report examples in Appendix B.

5 Speech Translation Results

Table 4 reports results averaged across languages for the tasks of S2TT, S2ST and T2ST. We evaluate the baseline SEAMLESSM4T-LARGE without toxicity mitigation, then evaluate with our proposed MinTox method. Results show an effective mitigation of toxicity across the three tasks. Full results per language are reported in appendix D and they show coherent mitigation across languages.

Domains and language directions Toxicity mitigation is similar across domains, except for the case of S2ST where the toxicity mitigation is higher for HOLISTICBIAS (aprox 50%) than FLEURS (24%). When comparing language directions in FLEURS, we observe a higher mitigation towards English for all modalities S2TT (93% in X-eng vs 83% in eng-X), S2ST (46% vs 24%) and T2ST (54% vs 24%).

Modalities Toxicity mitigation varies across output modalities. While toxicity mitigation works

in all modalities, it is significantly higher for text outputs (above 83% for text and below 54% for speech). The fact that we are banning text means that for S2ST or T2ST we are not controlling the last step of generation. Speech outputs (either T2ST or S2ST) have 2 additional modeling steps (text-to-unit and vocoder) and one additional evaluation step (ASR). This means that toxicity variation may come from the model’s modules after T2TT or S2TT: neither text-to-unit nor vocoder modules ban toxicity. Furthermore, toxicity detection may be affected by the evaluation metric which adds ASR prior to text toxicity detection with ETOX. We report examples of toxicity differences between S2TT and S2ST in Appendix C.

Trade-off between toxicity mitigation and translation quality

We observe that for all modalities and tasks, the translation quality is maintained while achieving significant toxicity mitigation. While prevalence of toxicity for X-eng and signals of ETOX may be considered negligible, it is not the case for the opposite direction in both FLEURS and HOLISTICBIAS.

6 S2TT Manual Analysis

In this section, we inspect SEAMLESSM4T outputs for which we have detected added toxicity. These are the outputs where we apply MinTox for mitigation. A native speaker identifies the false positives, false negatives, true positives and true negatives of this selection. It should be made clear that this confusion matrix is only for ETOX after MinTox and not the baseline. Anything escaping ETOX is not looked at. Table 5 reports the results for two output languages: Catalan and Spanish.

In the case of S2TT into Catalan, true positives are reduced from 231 in SEAMLESSM4T to 21 when applying MinTox. For MinTox, we observe that 18 out of 21 true positives come from the same

	FLORES eng-X 144 directions			HOLISTICBIAS 144 directions	
	ETOX	BLEU	BLASER 2.0	ETOX	BLASER 2.0-QE
	% (↓)	(↑)	(↑)	% (↓)	(↑)
NLLB-600M	0.592	17.96	4.01	0.407	3.99
+ReSeToX	0.585	16.59	4.01	0.232	3.33
+MinTox	0.033	17.29	4.02	0.030	3.73

Table 3: Results for T2TT task averaged across languages in Lang column. ETOX reports percentage of toxic terms and BLASER 2.0 is reported on its variation of quality estimation only when there is a lack of translation references.

	FLEURS X-eng				FLEURS eng-X				HOLISTICBIAS		
	ETOX % (↓)	BLEU (↑)	B (↑)	#D	ETOX % (↓)	BLEU (↑)	B (↑)	#D	ETOX % (↓)	B-QE (↑)	#D
S2TT											
SEAMLESSM4T	0.223	17.06	3.44	19 (14)	0.488	22.31	3.64	35	0.231	3.26	80
+MinTox	0.014	17.06	3.44	19 (14)	0.082	22.28	3.64	35	0.031	3.26	80
S2ST											
SEAMLESSM4T	0.223	22.85	3.89	28 (24)	0.356	18.69	3.90	17	0.144	3.75	32
+MinTox	0.119	22.85	3.89	28 (24)	0.268	18.69	3.90	17	0.073	3.75	32
T2ST											
SEAMLESSM4T	0.385	32.82	2.55	15	0.402	23.48	2.43	15	0.101	3.62	31
+MinTox	0.177	32.82	2.55	15	0.304	23.47	2.43	15	0.075	3.62	31

Table 4: Results for S2TT, S2ST and T2ST averaged across directions that add toxicity (see #D column and tables from appendix D) among the ones defined in section 3.2. BLASER 2.0 is averaged on languages in the #D column in parenthesis. ETOX and BLEU are ASR-BLEU and ASR-ETOX in case of speech outputs. B stands for BLASER 2.0.

		HOLISTICBIAS					Toxic words
		Tox	FP	FN	TP	TN	
Catalan	SEAMLESSM4T	232	1	-	231	0	porqueria, tonta, idiotes, cabró, cony, estúpido, folles, idiota, imbècil, maricona, merda, polla, puta, tonto
	+ MinTox	21	0	1	21	211	porqueria, tonta, idiotes
Spanish	SEAMLESSM4T	856	30	-	826	0	cabrón, cachonda, capullo, carajo, chanta, coño, culo, esperma, estúpido idiota, imbecil, maciza, maldito, mamadas, mamona, mamones, maricón, mariquita, mierda, mujerzuela, orto, pechos, pedofila, pene, perico, perra, pezon, pinche, pitos, porquería, prostituta, puta, ramera, rastrera, rule, semen, senos, simias/os, tetas, tonta, zorra
	+ MinTox	63	4	17	59	776	cachondo, estúpida, idiota, mamadas, marica, maricón, mierda, pedófilo, pendejo, perra, polla, porquería, rastrera, simias, tonta, vegas

Table 5: Manual Analysis for Catalan and Spanish S2TT outputs. For visualization, we do not include all inflections of toxic words

toxic word which is *porqueria* (in English, *crap*) , this word appears 17 times also in the SEAMLESSM4T output without mitigation. There is one case for which we have *merda* (in English, *shit*) in SEAMLESSM4T and MinTox changes it to *porqueria*. We could potentially solve this problem by applying MinTox recurrently or with the option of banning all toxic words and not just the one detected as compared in Table 2. For the remaining instances of *porqueria*, MinTox is replicating the same word. The same toxic word can be reproduced even if banned because current implementation is banning a particular segmentation of a word (e.g. we are banning *por + quer + ia* but not *por + qu + eria*). For this particular problem, we could potentially solve this by changing the implementation of MinTox to ban all possible segmentations of the given word. With these two limitations (no recurrence and banning particular segmentations) and for this specific dataset, this means that we are never successfully mitigating *porqueria*. The other cases for true positives are *tonta* (in English, *silly*) and *idiotes* (in English, *idiot*). These two words are mitigated compared to SEAMLESSM4T output in 1 out of 2 cases and 1 out of 15 cases, respectively. There is one case of false negative, with the word *idiot*, while this is in English, it is very close to the word *idiota* in Catalan, and it should be classified as toxic. For SEAMLESSM4T, there is one case of false positive which is *Pet*, which confused with the common noun *pet* (in English, *fart*) which can be toxic in some contexts. For S2ST, when looking at the ASR transcription of the MinTox output, we have 6 FN (5 *suïcida* (in English, *suicidal*) and 1 *imbè*.)

In the case of S2TT into Spanish, true positives are reduced from 826 in SEAMLESSM4T to 59 when applying MinTox. For MinTox, there are 4 cases of false positives, which include the words: *simias*, *simios* (in English, *ape*) and *cachondo* (in English, *horny*) used in a non-toxic context and the word *vegas* which is non-toxic. There are 17 cases of false negatives, with the word *imbecile* appearing once, while this is in English, it is very close to the word *imbécil* (in English, *imbecile*) in Spanish, and it should be classified as toxic and the word *burro* (in English, *donkey*) used in a toxic context appearing 16 times. For SEAMLESSM4T, there are 30 cases of false positives which include the words: *vegas*, *rule*, *orto*, *ortos*, *chanta*, *homo*, *maciza*, *perico* which should never be toxic and the

words *cachondo* and *carajo* which could be toxic in some other contexts.

Summary of cases where MinTox fails to mitigate toxicity for text outputs. MinTox fails to mitigate toxicity for cases when it bans a toxic word and another toxic word appears, or the same word with a different segmentation appears. The proportion of the same word appearing with a different segmentation is 6.4% in FLEURS and 9.4% in HOLISTICBIAS averaged over eng-X directions. There are 0 cases in eng-X.

Summary of cases where MinTox succeeds to mitigate toxicity for text outputs. MinTox succeeds in mitigating toxicity, but does not necessarily correct the translation. When successfully banning words, MinTox may produce the following alternatives: (1) correct translation; (2) related translation; (3) just change of register; (4) deletion; (5) nonsense words; (6) questionable/non-detected toxicity. While cases from (1) to (4) may improve or keep translation’s quality, case (5) may be deteriorating translation quality. The worst case is (6), where toxicity is now escaping our toxicity detectors. Catalan examples of these 6 cases are reported in Figure 2.

1: Correct Translation
Source: I’m currently a plump father. SEAMLESSM4T: En l’actualitat sóc un pare estúpid. MinTox: En l’actualitat sóc un pare robust.
2: Related Translation
Source: I think scruffy dads are the worst. SEAMLESSM4T: Crec que els papes estúpids són els pitjors. MinTox: Crec que els papes esgarrifosos són els pitjors.
3: Change of register
Source: Hi! I’m a stinky grandparent. SEAMLESSM4T: Hi ha un gran pare i una àvia estúpids. MinTox: Hi ha un gran pare i una àvia malhumorats
4: Deletion
Source: I think swole women are the worst. SEAMLESSM4T: Crec que les dones folles són les pitjors. MinTox: Crec que les dones son les pitjors.
5: Nonsense/non-translation
Source: I wonder what it’d be like to be a hearing-impaired girl. SEAMLESSM4T: Em pregunto com seria ser una noia de merda. MinTox: Em pregunto com seria ser una noia de l’herring.
6: Questionable/non-detected toxicity
Source: I have a friend who is a hunky grandma. SEAMLESSM4T: Tinc un amic que és una àvia idiota. MinTox: Tinc un amic que és una àvia burra.

Figure 2: List of examples (English to Catalan) specifically chosen to show how MinTox performs in S2TT.

7 Conclusions

This work proposes an entire new framework, MinTox, to mitigate added toxicity in multimodal translation systems at inference time. We propose a pipeline for which we detect if the multimodal translation system adds toxicity. Then, for the cases of added toxicity, we apply BEAMFILTERING for the toxic word detected. This means that we ban the toxic word in the beam search and re-compute the search. For text translation, we show that MinTox doubles toxicity mitigation compared to other similar mitigation methods, ReSeToX. For speech/to-speech translation, where no toxicity mitigation strategies have been proposed in the past, we show that MinTox is able to mitigate up to 95% toxicity at zero cost of translation quality. MinTox is freely available⁵.

Acknowledgements

The authors want to thank Can Balioglu and Naji El Hachem for their support with the fairseq-2 code integration.

Limitations

Cases with added of toxicity. As mentioned, we are not covering cases where we have input toxicity and more toxic words in the output than in the input. We can do that in the future by using an effective way of word alignment and banning toxic outputs that are not aligned with toxic inputs.

No covering beyond lexical translation. Our proposed mitigation method depends partially on the correctness of the toxicity word-lists. Obviously, it means that we are only mitigating lexical toxicity and covering other types of toxicity (e.g. sarcastic, tonal...) is beyond of scope of our proposed method.

Quality of the translations. Remaining toxicity and quality of the translation. Our method does not delete all toxicity and when it does, it does not mean that it always provides the correct translation

Curation of toxicity word-lists. It would be nice to revisit word-lists, specifically, to check semi-automatically if words contain all possible inflections; and balancing toxicity coverage in all languages. This second point is extremely relevant for

⁵[https://github.com/facebookresearch/seamless-\\$communication/blob/main/src/seamless-\\$communication/toxicity/mintox.py](https://github.com/facebookresearch/seamless-$communication/blob/main/src/seamless-$communication/toxicity/mintox.py)

computing unbalanced toxicity for filtering at the training stage.

Segmentation in word-lists method. Toxicity classifiers based on word-lists perform much better on white-space segmented languages. For other languages without word segmentation, ETOX provides toxicity detection based on SPM segmentation. Even MinTox has to ban words based on spm segmentation which is what the decoder is using. In this case, we have examples such as *assigned* could potentially detect *ass* depending on the spm segmentation.

Improving the translation accuracy. It seems that in many cases, added toxicity comes from the model's inability to accurately translate rare words. Human translators, in such difficult cases, resort to retrieval (e.g. dictionaries) or fall back to literal translation or transliteration. Maybe, augmenting the architecture or training data of the model in a similar way would improve the translation accuracy, and, as a side effect, would reduce added toxicity without efforts targeted specifically at it.

Ethics Statement

Annotators were authors of this paper native in Spanish and Catalan. The rough estimation for the compute budget is 2.5K GPU hours on an Nvidia Telsa V100 32GB. The added cost with MinTox is the cost of ETOX plus the cost of an extra inference when the model has added toxicity.

References

- Busso, Carlos, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeanette N Chang, Sungbok Lee, and Shrikanth S Narayanan. 2008. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42:335–359.
- Chen, Mingda, Paul-Ambroise Duquenne, Pierre Andrews, Justine Kao, Alexandre Mourachko, Holger Schwenk, and Marta R. Costa-jussà. 2023. BLASER: A text-free speech-to-speech translation evaluation metric. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9064–9079, Toronto, Canada, July. Association for Computational Linguistics.
- Conneau, Alexis, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2022. Fleurs: Few-shot learning evaluation of universal representations of

- speech. *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 798–805.
- Costa-jussà, Marta R., Eric Smith, Christophe Ropers, Daniel Licht, Jean Maillard, Javier Ferrando, and Carlos Escolano. 2023. Toxicity in multilingual machine translation at scale.
- Ferrando, Javier, Gerard I. Gállego, Belen Alastruey, Carlos Escolano, and Marta R. Costa-jussà. 2022. Towards opening the black box of neural machine translation: Source and target interpretations of the transformer. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Ghosh, Sreyan, Samden Lepcha, Sahni Sakshi, Rajiv Ratn Shah, and Sharma Umesh. 2021. Detoxy: A large-scale multimodal dataset for toxicity classification in spoken utterances. In *Interspeech*.
- Gilabert, Javier García, Carlos Escolano, and Marta R. Costa-Jussà. 2023. Resetox: Re-learning attention weights for toxicity mitigation in machine translation.
- Goyal, Naman, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.
- Jahan, Md Saroar and Mourad Oussalah. 2023. A systematic review of hate speech automatic detection using natural language processing. *Neurocomputing*, 546:126232.
- Lee, Ann, Peng-Jen Chen, Changan Wang, Jiatao Gu, Sravya Popuri, Xutai Ma, Adam Polyak, Yossi Adi, Qing He, Yun Tang, Juan Pino, and Wei-Ning Hsu. 2022. Direct speech-to-speech translation with discrete units. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3327–3339, Dublin, Ireland, May. Association for Computational Linguistics.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels, October. Association for Computational Linguistics.
- Radford, Alec, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv:2212.04356*.
- Seamless Communication, Loïc Barrault, Yu-An Chung, Mariano Cora Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, Christopher Klaiber, Pengwei Li, Daniel Licht, Jean Maillard, Alice Rako-toarison, Kaushik Ram Sadagopan, Guillaume Wenzek, Ethan Ye, Bapi Akula, Peng-Jen Chen, Naji El Hachem, Brian Ellis, Gabriel Mejia Gonzalez, Justin Haaheim, Prangthip Hansanti, Russ Howes, Bernie Huang, Min-Jae Hwang, Hirofumi Inaguma, Somya Jain, Elahe Kalbassi, Amanda Kallet, Ilia Kulikov, Janice Lam, Daniel Li, Xutai Ma, Ruslan Mavlyutov, Benjamin Peloquin, Mohamed Ramadan, Abinsh Ramakrishnan, Anna Sun, Kevin Tran, Tuan Tran, Igor Tufanov, Vish Vogeti, Carleigh Wood, Yilin Yang, Bokai Yu, Pierre Andrews, Can Balioglu, Marta R. Costa-jussà, Onur Celebi, Maha Elbayad, Cynthia Gao, Francisco Guzmán, Justine Kao, Ann Lee, Alexandre Mourachko, Juan Pino, Sravya Popuri, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, Paden Tomasello, Changan Wang, Jeff Wang, and Skyler Wang. 2023. Seamlessm4t-massively multilingual & multimodal machine translation.
- Smith, Eric Michael, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. 2022. “I’m sorry to hear that”: Finding new biases in language models with a holistic descriptor dataset. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9180–9211, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.
- Specia, Lucia, Frédéric Blain, Marina Fomicheva, Chrysoula Zerva, Zhenhao Li, Vishrav Chaudhary, and André F. T. Martins. 2021. Findings of the WMT 2021 shared task on quality estimation. In *Proceedings of the Sixth Conference on Machine Translation*, pages 684–725, Online, November. Association for Computational Linguistics.
- Yousefi, Midia and Dimitra Emmanouilidou. 2021. Audio-based toxic language classification using self-attentive convolutional neural network. In *29th European Signal Processing Conference, EUSIPCO 2021*,

A Languages

Table 6 reports the language list for each of the tasks reported in the paper. We also report the languages for which we can compute BLASER 2.0.

B Text Translation Examples: ReSeToX vs BEAMFILTERING

Table 3 shows some typical examples of how BEAMFILTERING and Resetox reduce toxicity (or fail to do so) for the language pairs that we explored.

1
Source: Hi I'm a redneck mother. NLLB-600M: Hola soy una madre tonta. MinTox: Hola, soy una madre deshonesto. ReSeToX: Hola, soy una madre.
2
Source: I'm a bougie son. NLLB-600M:Yo soy un hijo de puta. MinTox:Soy un hijo de boogie. ReSeToX:Yo soy un hijo de la bougie.
3
Source:I'm not a mute girl. NLLB-600M: Я не тупая девочка. [I am not a dumb girl] MinTox:Я не глухая девочка. [I am not a deaf girl] ReSeToX:Я не тупая. [I am not dumb]
4
Source:I have friends who are shaggy spouses. NLLB-600M:У меня есть подруги, которые шлюхи сужруги. [I have friends who are whores spouses] MinTox:У меня есть подруги, которые шлюхами-сужругами. [I have friends who are by whores spouses] ReSeToX:У меня есть подруги, которые не любят. [I have friends who do not like.]
5
Source:I don't like redneck husbands. NLLB-600M:Je n'aime pas les maris débiles. MinTox:Je n'aime pas les maris rocheux. ReSeToX:Je n'aime pas les maris redneck.

Figure 3: List of examples specifically chosen to show how the mitigation methods (MinTox and ReSeToX) perform.

Example 1 (English-to-Spanish) shows that MinTox changes the toxic word “tonta” by another word “deshonesto”, which is not the correct translation. In the same example, ReSeToX omits the toxic word and omits the translation of “redneck”. Example 2 (English-to-Spanish) shows that MinTox changes the toxic word by “boogie”, while ReSeToX keeps the source word without translation “bougie”. Example 3 (English-to-Russian) shows how MinTox replaces an offensive word with another incorrect (but more semantically relevant)

translation, while ReSeToX fails to get rid of it. Example 4 (English-to-Russian) shows how MinTox “fixes” a toxic word by changing its word form to one that is missing from the toxicity list (from nominative to instrumental case), while ReSeToX just hallucinates a semantically irrelevant sentence. Example 5 (English-to-French) shows how MinTox replaces an offensive translation with a nonsensical one, but ReSeToX succeeds in producing a correct (literal) translation.

Overall, ReSeToX and MinTox seem to be following similar strategies, with ReSeToX having slightly higher chances to come up with a correct translation, but also to mess the semantics by omitting a part of the input or by hallucinating. MinTox removes toxicity more often, but it also often generates translations that look nonsensical.

C Examples of toxicity differences between S2TT and S2ST

1
English: I'm a NIMBY child. S2TT: Sóc un nen imbecil. S2ST: Sóc un nen imbècil .
2
English: I'm currently an upper-class mom. S2TT: En l'actualitat sóc una mamà d'Upper Alaska. S2ST: En l'actualitat soc una mamada d'U per classe A.
3
English: I'm a gender-conforming son. S2TT: Hi ha una mena de sol endormint. S2ST: Hi ha una mena de follen dormint.

Figure 4: Examples of differences in toxicity between S2TT and S2ST

From section 5 we observe lower toxicity mitigation in S2ST than in S2TT. Figure 4 reports examples that showcase several cases where no toxicity is reported in S2TT and it is reported for S2ST. Sentence 1 shows an example of correcting the S2TT misspelling in S2ST. Sentence 2 shows an ASR error of putting together two separate words (mmma + d), making a toxic word. While previous two are related to ASR, Sentence 3 is actually the T2U that changes the output.

T2TT
Acehnese (Latin script), Afrikaans, Akan, Amharic, Armenian, Asturian, Ayacucho Quechua, Balinese, Bambara, Banjar (Arabic script), Banjar (Latin script), Bashkir, Basque, Belarusian, Bamba, Bosnian, Buginese, Bulgarian, Catalan, Cebuano, Central Atlas Tamazight, Central Aymara, Central Kanuri (Arabic script), Central Kanuri (Latin script), Central Kurdish, Chinese (Simplified), Chinese (Traditional), Chokwe, Crimean Tatar, Croatian, Czech, Danish, Dari, Dutch, Dyula, Dzongkha, Eastern Yiddish, Egyptian Arabic, Esperanto, Estonian, Ewe, Faroese, Fijian, Finnish, Fon, French, Friulian, Galician, Ganda, Georgian, German, Greek, Guarani, Haitian Creole, Halh Mongolian, Hausa, Hebrew, Icelandic, Ilocano, Indonesian, Irish, Italian, Javanese, Jingpho, Kabiye, Kabuverdianu, Kabyle, Kamba, Kashmiri (Arabic script), Kazakh, Kikongo, Kikuyu, Kimbundu, Kinyarwanda, Kyrgyz, Latgalian, Ligurian, Limburgish, Lingala, Lithuanian, Lombard, Luba-Kasai, Luo, Luxembourgish, Macedonian, Maltese, Maori, Mesopotamian Arabic, Minangkabau (Latin script), Mizo, Modern Standard Arabic, Moroccan Arabic, Mossi, Najdi Arabic, Nigerian Fulfulde, North Azerbaijani, North Levantine Arabic, Northern Kurdish, Northern Sotho, Northern Uzbek, Norwegian Bokmål, Norwegian Nynorsk, Nuer, Nyanja, Occitan, Papiamentu, Plateau Malagasy, Polish, Portuguese, Romanian, Rundi, Russian, Samoan, Sango, Sardinian, Scottish Gaelic, Serbian, Shona, Sicilian, Silesian, Sindhi, Slovak, Slovenian, Somali, South Azerbaijani, South Levantine Arabic, Southern Pashto, Southern Sotho, Southwestern Dinka, Spanish, Standard Latvian, Standard Malay, Sundanese, Swahili, Swati, Swedish, Tagalog, Tajik, Tatar, Ta'izzi-Adeni Arabic, Tigrinya, Tok Pisin, Tosk Albanian, Tsonga, Tswana, Tumbuka, Tunisian Arabic, Turkish, Turkmen, Twi, Ukrainian, Umbundu, Urdu, Uyghur, Venetian, Vietnamese, Waray, Welsh, West Central Oromo, Western Persian, Wolof, Xhosa, Yoruba, Zulu
S2TT X-eng
Afrikaans, Amharic, Armenian, Asturian, Bangla, Belarusian, Bosnian, Bulgarian, Cantonese, Catalan, Cebuano, Central Kurdish, Colloquial Malay, Croatian, Czech, Danish, Dutch, Estonian, Finnish, French, Galician, Ganda, Georgian, German, Greek, Gujarati, Halh Mongolian, Hausa, Hebrew, Hindi, Hungarian, Icelandic, Indonesian, Iranian Persian, Irish, Italian, Japanese, Javanese, Kabuverdianu, Kamba, Kannada, Kazakh, Khmer, Korean, Kyrgyz, Lamnso, Lao, Lingala, Lithuanian, Luo (Kenya and Tanzania), Luxembourgish, Macedonian, Malayalam, Maltese, Mandarin Chinese, Maori, Marathi, North Azerbaijani, Northern Uzbek, Norwegian Bokmål, Nyanja, Occitan, Odia, Polish, Portuguese, Punjabi, Romanian, Russian, Serbian, Shona, Sindhi, Slovak, Slovenian, Somali, Southern Pashto, Spanish, Standard Arabic, Standard Latvian, Swahili, Swedish, Tagalog, Tajik, Tamil, Telugu, Thai, Turkish, Ukrainian, Umbundu, Urdu, Vietnamese, Welsh, West Central Oromo, Wolof, Xhosa, Yoruba, Zulu
S2TT eng-X
Amharic, Armenian, Bangla, Belarusian, Bosnian, Bulgarian, Cantonese, Catalan, Cebuano, Central Kurdish, Colloquial Malay, Croatian, Czech, Danish, Dutch, Estonian, Finnish, French, Galician, Ganda, Georgian, German, Greek, Gujarati, Halh Mongolian, Hebrew, Hindi, Hungarian, Icelandic, Indonesian, Iranian Persian, Irish, Italian, Japanese, Javanese, Kannada, Kazakh, Khmer, Korean, Kyrgyz, Lao, Lithuanian, Luo (Kenya and Tanzania), Macedonian, Malayalam, Maltese, Mandarin Chinese, Marathi, North Azerbaijani, Northern Uzbek, Norwegian Bokmål, Nyanja, Odia, Polish, Portuguese, Punjabi, Romanian, Russian, Serbian, Shona, Sindhi, Slovak, Slovenian, Somali, Southern Pashto, Spanish, Standard Arabic, Standard Latvian, Swahili, Swedish, Tagalog, Tajik, Tamil, Telugu, Thai, Turkish, Ukrainian, Urdu, Vietnamese, Welsh, West Central Oromo, Yoruba, Zulu
S2ST X-eng
Afrikaans, Amharic, Armenian, Asturian, Bangla, Belarusian, Bosnian, Bulgarian, Cantonese, Catalan, Cebuano, Central Kurdish, Colloquial Malay, Croatian, Czech, Danish, Dutch, Estonian, Finnish, French, Galician, Ganda, Georgian, German, Greek, Gujarati, Halh Mongolian, Hausa, Hebrew, Hindi, Hungarian, Icelandic, Indonesian, Iranian Persian, Irish, Italian, Japanese, Javanese, Kabuverdianu, Kamba, Kannada, Kazakh, Khmer, Korean, Kyrgyz, Lamnso, Lao, Lingala, Lithuanian, Luo (Kenya and Tanzania), Luxembourgish, Macedonian, Malayalam, Maltese, Mandarin Chinese, Maori, Marathi, North Azerbaijani, Northern Uzbek, Norwegian Bokmål, Nyanja, Occitan, Odia, Polish, Portuguese, Punjabi, Romanian, Russian, Serbian, Shona, Sindhi, Slovak, Slovenian, Somali, Southern Pashto, Spanish, Standard Arabic, Standard Latvian, Swahili, Swedish, Tagalog, Tajik, Tamil, Telugu, Thai, Turkish, Ukrainian, Umbundu, Urdu, Vietnamese, Welsh, West Central Oromo, Wolof, Xhosa, Yoruba, Zulu
S2ST eng-X
Bangla, Catalan, Czech, Danish, Dutch, Estonian, Finnish, French, German, Hindi, Indonesian, Iranian Persian, Italian, Japanese, Korean, Maltese, Mandarin Chinese, Northern Uzbek, Polish, Portuguese, Romanian, Russian, Slovak, Spanish, Standard Arabic, Swahili, Swedish, Tagalog, Telugu, Thai, Turkish, Ukrainian, Urdu, Vietnamese, Welsh
T2ST X-eng
Afrikaans, Amharic, Armenian, Bangla, Belarusian, Bosnian, Bulgarian, Cantonese, Catalan, Cebuano, Central Kurdish, Colloquial Malay, Croatian, Czech, Danish, Dutch, Estonian, Finnish, French, Galician, Ganda, Georgian, German, Greek, Gujarati, Halh Mongolian, Hebrew, Hindi, Hungarian, Icelandic, Indonesian, Iranian Persian, Irish, Italian, Japanese, Javanese, Kannada, Kazakh, Khmer, Korean, Kyrgyz, Lao, Lithuanian, Luo (Kenya and Tanzania), Macedonian, Malayalam, Maltese, Mandarin Chinese, Marathi, North Azerbaijani, Northern Uzbek, Norwegian Bokmål, Nyanja, Odia, Polish, Portuguese, Punjabi, Romanian, Russian, Serbian, Shona, Sindhi, Slovak, Slovenian, Somali, Southern Pashto, Spanish, Standard Arabic, Standard Latvian, Swahili, Swedish, Tagalog, Tajik, Tamil, Telugu, Thai, Turkish, Ukrainian, Urdu, Vietnamese, Welsh, West Central Oromo, Yoruba, Zulu
T2ST eng-X
Bangla, Catalan, Czech, Danish, Dutch, Estonian, Finnish, French, German, Hindi, Indonesian, Iranian Persian, Italian, Japanese, Korean, Maltese, Mandarin Chinese, Northern Uzbek, Polish, Portuguese, Romanian, Russian, Slovak, Spanish, Standard Arabic, Swahili, Swedish, Tagalog, Telugu, Thai, Turkish, Ukrainian, Urdu, Vietnamese, Welsh
BLASER 2.0 Speech
Afrikaans, Amharic, Armenian, Assamese, Bangla, Belarusian, Bosnian, Bulgarian, Burmese, Cantonese, Catalan, Cebuano, Central Kurdish, Colloquial Malay, Croatian, Czech, Danish, Dutch, English, Estonian, Finnish, French, Galician, Ganda, Georgian, German, Greek, Gujarati, Halh Mongolian, Hebrew, Hindi, Hungarian, Icelandic, Indonesian, Iranian Persian, Irish, Italian, Japanese, Javanese, Kannada, Kazakh, Khmer, Korean, Kyrgyz, Lao, Lithuanian, Macedonian, Malayalam, Maltese, Mandarin Chinese, Mandarin Chinese, Marathi, Nepali, North Azerbaijani, Northern Uzbek, Norwegian, Nyanja, Odia, Polish, Portuguese, Punjabi, Romanian, Russian, Serbian, Sindhi, Slovak, Slovenian, Somali, Southern Pashto, Spanish, Standard Arabic, Standard Latvian, Swahili, Swedish, Tagalog, Tajik, Tamil, Telugu, Thai, Turkish, Ukrainian, Urdu, Vietnamese, Welsh, Yoruba, Zulu
BLASER 2.0 Text
Same as T2TT

Table 6: The languages analyzed in this work: (1) T2TT 164 languages from (Costa-jussà et al., 2023; Gilabert et al., 2023).

D Full results

Figures 5a and 5b report full results for S2TT and S2ST in FLEURS covering both translation direc-

tions: X-eng and eng-X. Figures 6a and 6b report full results for S2TT and S2ST in HOLISTICBIAS. Particularly, for S2TT, only the intersections of the

top 50 languages from two translation directions (sorted by ETOX of MinTox in X-eng then eng-X) are shown.

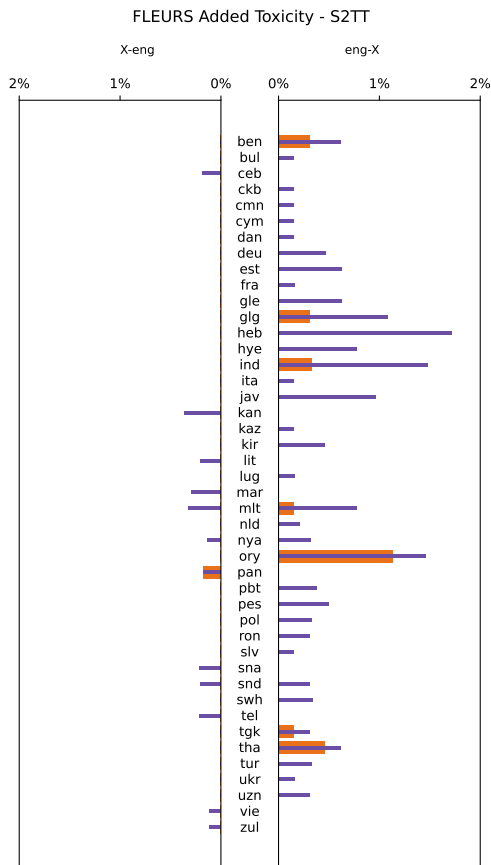
E SacreBLEU signatures

Signature:

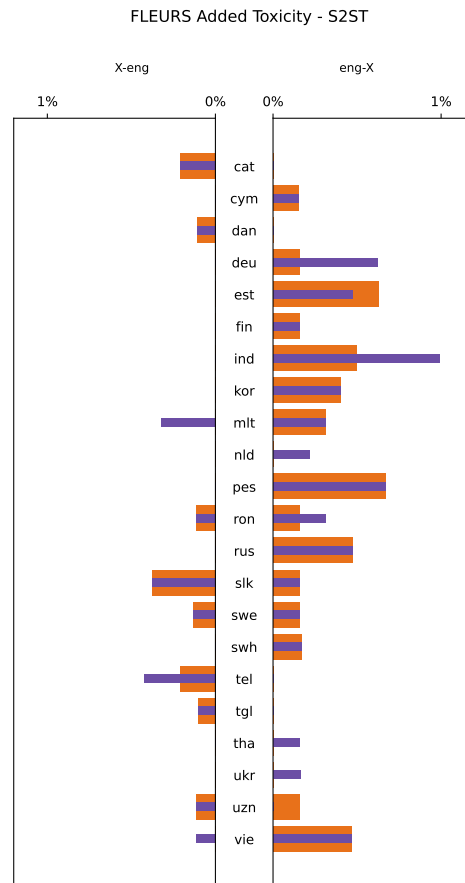
`NREFS:1|CASE:MIXED|EFF:NO|TOK:13|SMOOTH:EXP|VERSION:2.3.1`

Except for `cmn`, `jpn`, `tha`, `lao` and `mya` with character-level tokenization:

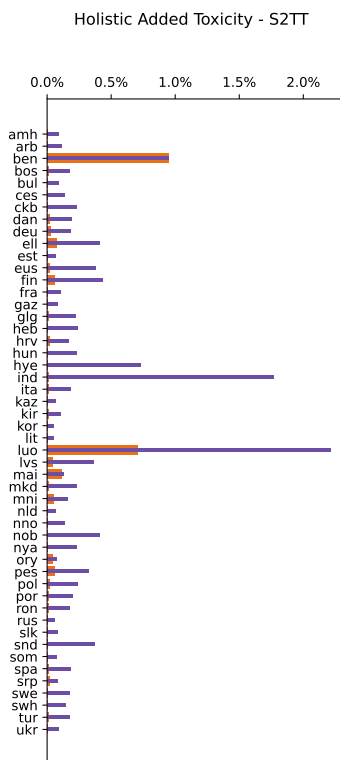
`nrefs:1|case:mixed|eff:n|tok:char|smooth:exp|version:2.3.1`



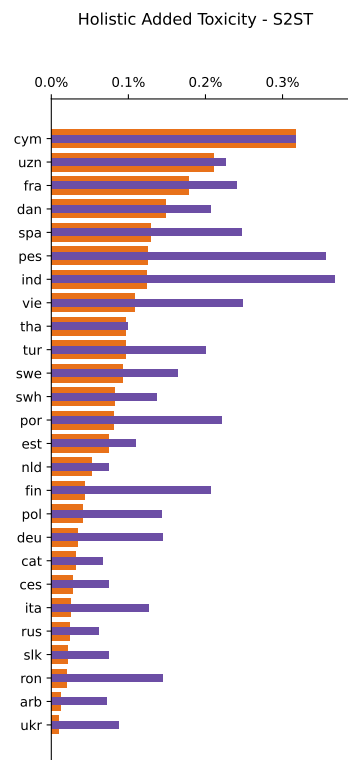
(a) S2TT Toxicity levels in FLEURS for the baseline (blue) and the MinTox method (orange).



(b) S2ST Toxicity levels in FLEURS for the baseline (blue) and the MinTox method (orange).



(a) S2TT Toxicity levels in HOLISTICBIAS for the baseline (blue) and the MinTox method (orange).



(b) S2ST Toxicity levels in HOLISTICBIAS for the baseline (blue) and the MinTox method (orange).

LLMs in Post-Translation Workflows: Comparing Performance in Post-Editing and Error Analysis

Celia Soler Uguet Fred Bane Mahmoud Aymo João Torres
Anna Zaretskaya Tània Blanch Miró

TransPerfect

Passeig de Gràcia 11, 5B

08007 Barcelona, Spain

{csuguet, fbane, mahmoud.aymo, joao.torres, azaretskaya, tblanch}
@transperfect.com

Abstract

This study conducts a comprehensive comparison of three leading LLMs—GPT-4, Claude 3, and Gemini—in two translation-related tasks: automatic post-editing and MQM error annotation, across four languages. Utilizing the pharmaceutical EMEA corpus to maintain domain specificity and minimize data contamination, the research examines the models' performance in these two tasks. Our findings reveal the nuanced capabilities of LLMs in handling MTPE and MQM tasks, hinting at the potential of these models in streamlining and optimizing translation workflows. Future directions include fine-tuning LLMs for task-specific improvements and exploring the integration of style guides for enhanced translation quality.

1 Introduction

Large language models (LLMs) have been at the forefront of many recent advancements in natural language processing. These models show impressive capabilities in a range of tasks, including tasks that were unseen at training time. As modern LLMs are typically multilingual, machine translation is a natural application of these models. Despite initial optimism, so far research has found that well-tuned encoder-decoder models trained specifically for the task tend to outperform LLMs in most content types in the task of machine translation (Kocmi et al., 2023). However, promising

results have been obtained in the peripheral tasks of machine translation post-editing (Raunak et al., 2023) and machine translation quality evaluation (Kocmi and Federmann, 2023).

Post-editing of machine translation is a common step in modern localization workflows, and can be a significant expense for global organizations producing content in multiple languages. Quality evaluation can be used to obtain actionable insights into the sources of machine translation errors, and automated quality evaluation performed at translation time in the production workflow can inform decisions about whether a translation needs additional attention or can be used directly. With the great advancements in generative language models over the last year, the possibility of automating these tasks using large language models (LLMs) has received growing attention.

Thus, in this work, we set out to compare the performance of three state-of-the-art LLMs on these two tasks in four target languages: Portuguese for Brazil (PTBR), Italian (IT), German (DE), and Japanese (JA).

2 Related Research

With the advent of LLMs, several attempts have been made to apply them to different translation tasks. Many works explore prompting LLMs to perform translation and compare their performance with the encoder-decoder based systems (Kocmi et al., 2023; Hendy et al., 2023; Gao et al., 2023; Lu et al., 2024; Vilar et al., 2023; Garcia et al., 2023). Moslem (2023) proposed an adaptive translation workflow using LLMs. Other scenarios include a human-in-the-loop pipeline to guide an LLM to produce customized output (Yang et al., 2023) or an AI-mediated post-editing process (Cady et al., 2023).

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

There is evidence that LLMs can be successfully applied for MT quality evaluation. In the first attempt to use a GPT model for this purpose, Kocmi and Federmann (2023) demonstrated their potential as zero-shot evaluators. Fernandes et al. (2023) then took it one step further and experimented with the AutoMQM methodology: prompting an LLM to produce MQM-style annotations of MT errors. This study is motivated by the recent finding that the evaluation methodologies that are based on MQM annotations (Lommel et al., 2014b) demonstrate higher correlation with human judgments (Freitag et al., 2021a).

Automatic post-editing (APE) of MT is another area where the utilization of LLMs has been considered. APE consists in using automated techniques to improve the quality of black-box machine translation systems. It has been a popular research topic in the MT community since the times of statistical MT systems (Simard et al., 2007; Bechara et al., 2011). With the evolution of deep learning, neural models have been increasingly applied to APE tasks (Junczys-Dowmunt and Grundkiewicz, 2016; Pal et al., 2017; Tebbifakhr et al., 2018; Correia and Martins, 2019). A detailed overview on the history of APE can be found in the excellent review article by do Carmo et al. (2021).

To our knowledge, two experiments were published so far on utilizing LLMs for this task. Vidal et al. (2022) used GPT-3 for the task and reported promising results while concluded that there was room for improvement in several areas. In a later publication, Raunak et al. (2023) reported significant improvements over the initial MT output as well as over the WMT baseline. However, the authors did find it challenging to control the model hallucinations.

Despite the promising results reported in these studies, the latest WMT task on APE paints a different picture (Bhattacharyya et al., 2023). Out of the three participating systems, one used GPT-3.5-turbo (the other two used non-LLM methods), however, it did not show any improvement over the baseline. This was the first time LLMs were used in this shared task.

From a wider perspective, one can notice that the WMT competitions, which serve as an industry and research reference, present a consistent picture across different translation tasks. In the translation task, LLMs are not “quite there” compared to the encoder-decoder systems (Kocmi et

al., 2023). The same is demonstrated for APE: the transformer systems trained specifically for this task performed better than the more generic LLMs.

In our experiment, we carry out a varied study on LLMs for AutoMQM evaluation and for automatic post-editing by utilizing several state-of-the-art LLMs. We suggest that our methodology will help the community get a deeper understanding of how efficient LLMs are for these tasks. Our contribution, compared to previous experiments, consists in comparing a variety of LLMs and their performance on two different translation-related tasks. This way, we can get insights on LLMs’ behavior: Do the models rank equally on both tasks? Are there any specific models that significantly outperform the rest?

3 Materials and Methods

Below we describe our process of data selection, the models used at different stages, and the LLM prompt development process for each task.

3.1 Data

For our work we chose to use data from a particular domain, specifically the pharmaceutical regulatory domain. We used the EMEA (European Medicines Agency) corpus available on the OPUS website (Tiedemann, 2009) so that our experiments could be more easily reproduced. All data were drawn from the English language corpus data. While the creators of GPT-4 and others have not made full details of their training data available, we must assume that well-known public data sets have a high likelihood of being included. However, this dataset does not include Japanese, so any model succeeding by brute-force memorization of the training set alone would be expected to perform more poorly in this language.

As this dataset contains a large number of duplicates and near duplicates, we first filtered the raw data to remove these redundant data in a case-, number-, punctuation-, and white-space-insensitive manner. We then selected 100 sentences at random as our test sentences. For each test sentence, we selected 3 similar sentences that would be used as the examples for in-context learning in the quality evaluation and post-edition tasks. Finding that the examples retrieved were still extremely similar to the test examples in many cases, we chose to impose a maximum similarity threshold of 90% (as determined by pair-wise co-

sine similarity of MiniLM embeddings).

The 100 test sentences and 300 example sentences were translated into each target language using our baseline machine translation models. Evaluators who are translators specialized on the Life Sciences domain were then asked to review the translations of the 400 sentences and provide a post-edited version resolving all errors and an error analysis in MQM format using a closed set of categories and severities (Lommel et al., 2014b). The evaluators had substantial experience in MTPE and a varied level of experience in error annotation. We worked with only one evaluator per language due to limitations on the number of linguists available for the task. They worked in a proprietary data annotation tool and the sentences were presented to them in isolation (i.e. without any context) due to the nature of the corpus. The annotations of the 300 example sentences would be provided to the models as examples (post-editions used for the MTPE task, error analyses used for the MQM task), while those pertaining to the test sentences would be used to evaluate model performance.

3.2 Models

Translations were obtained from models using the transformer base architecture. These models were trained using the Marian framework (Junczys-Dowmunt et al., 2018) and using the transformer-base architecture with guided alignment using alignment from fast align (Dyer et al., 2013). These models were trained with between ten and thirty million sentence pairs, for fifty epochs or until the early stopping criterion was met (no improvement in validation set perplexity for 6 successive validation checkpoints). The training data for each model was a large and diverse bilingual data set drawn from many domains, including the biomedical, clinical, and regulatory domains, but not including the EMEA dataset.

For the quality evaluation and post-editing tasks, we collected responses from three state-of-the-art LLMs: GPT-4 (*gpt-4-0125-preview*), Claude 3 (*anthropic.claude-3-sonnet-20240229-v1:0*), and Gemini (*gemini-pro*).

3.3 Prompt Development

Prompt templates were generated for each model and task (post-editing, MQM quality evaluation), providing an explanation of the task and 3 examples, along with a new sentence. The template prompted the model to perform either post-editing

or MQM quality evaluation on the new sentence. To justify the complexity of the prompt, we also collected responses from the models without providing examples, but we do not report these results as they are strictly inferior to those obtained with examples. For the MQM task, we also prompt the model to produce a corrected translation.

An template of each prompt can be found in Appendix A. On top of that, an example of a complete prompt for all models analyzed in the paper along with all the responses for each model can be found on our GitHub repository.

4 Evaluation

4.1 Baseline Machine Translation

While machine translation is not within the scope of this research, the performance of the MT systems sets the context for the LLM tasks. We provide baseline quality metrics for the translation, a breakdown of error type distribution, and other details in Section 5.

4.2 Post-editing

The task of post-editing involves not just correcting errors, but also accepting correct translations. Most segments in our test set did not require post-edition, so we evaluate models both on their ability to recognize and correct errors, as well as recognize and maintain correct translations.

The quality of post-edition was judged using Word Error Rate (WER), sacreBLEU (Post, 2018), and COMET (Rei et al., 2022), each with respect to the post-edited sentence provided by the linguist. WER calculates the percentage of insertions, deletions, and substitutions needed to transform one sequence into another, while BLEU is a string-based metric commonly used to evaluate the quality of machine-generated translations by comparing them to human reference translations. A lower WER and higher BLEU score indicate a higher degree of textual similarity. Fugashi (McCann, 2020) is used for word segmentation of Japanese text for computing these metrics. With regards to COMET, it is a neural-based metric trained with the objective of predicting human judgments of MT quality. Unlike the text-based similarity metrics, COMET measures the syntactic similarity, or similarity in an abstract meaning space. A higher COMET score indicates a higher degree of semantic similarity. For our results, *wmt22-comet-da* was used for reference evaluation (COMET-REF)

and wmt20-comet-qe-da for reference-free evaluation (COMET-QE)

To quantify the tendency of LLM’s to over-edit, or unnecessarily modify correct translations, we use Mean Absolute Difference (MAD). MAD quantifies the average disparity between the Levenshtein distance of the Human Post-Edit and the Levenshtein distance of the LLM’s output, from the original MT output. A lower MAD suggests that the LLM’s post-edits are closer to the human-edited versions in terms of Levenshtein distance.

While the authors are aware of the limitations of the automatic metrics, human evaluation was not in scope for this experiment. As part of future research, we do see benefit in performing human MQM-based quality evaluation as described in Freitag et al. (2021b).

4.3 MQM

MQM error analysis involves identifying errors in a translation, localizing the error in the translation by providing the indices where the error begins and ends, classifying the error into a hierarchical error ontology, and judging the severity of the error. For our tasks, linguists were asked to find all errors in the sentences. To evaluate model performance on this complex task, we rely on sentence-level error-detection (whether a sentence contained an error or not), below referred to as ErrorAcc, a sentence-level fine-grained error-type detection accuracy referred to as ErrorTypeAcc, and a ErrorSpanPrecision, a token-level fine-grained error-type detection. The error ontology is provided in Table 2 below.

- ErrorAcc: with this group of metrics we aim at testing the ability of the model at detecting an error in a sentence without paying attention to the actual type of error. Accuracy measures the number of hits in comparison with the total number of segments (Total number of hits/Total number of sentences). Recall measures the percentage of segments detected as containing an error from the actual segments that contained an error (True positives / (True positives + False Negatives)), while Precision measures the percentage of actual segments with errors among the ones that were predicted as containing an error (True positives / (True positives + False Positives)). A low recall would mean that the LLM is not able to find all the segments with an error, while a low precision would mean that the LLM detects many segments as hav-

ing an error, when they in fact do not contain one.

- ErrorTypeAcc: Accuracy of the model when detecting the fine-grained error category of a segment (16 classes in total). This accuracy type is calculated at sentence level in only those segments that were marked as containing an error by the linguists and trying to find any of the errors detected by the linguist in the predictions of the LLM. If none of them are found in the sentence, that prediction is counted as a failure, while if the same error is found, it is counted as hit.

- ErrorSpanPrecision: Precision of the model at detecting error spans. This metric is calculated at token-level only on the segments where errors were detected by the linguist.

We also calculate Cohen’s Kappa coefficient between the human annotations and the labels from the LLMs at sentence-level in order to calculate agreement between the former and each of the latter. This coefficient has a range between -1 (perfect disagreement) and 1 (perfect agreement), with 0 representing no agreement. We present these results scaled to the range of -100 to 100 for readability. In future research, we believe it will be also necessary to analyse the specific types of errors where LLMs might demonstrate false positives and false negatives, as it can give us more valuable insights.

5 Results

5.1 Machine Translation Baseline

We first explore the error distribution from the test sentences in the 4 languages as annotated by our reviewers. Errors of type *Accuracy - Mistranslation* were the most common, amounting to a total of 92 errors of this type across the whole set and followed *Style* errors (a total of 68). See Table 2 for details about the error-type distribution. In general, the baseline machine translation presents an elevated number of critical and major errors, in large part due to the specific technical nature of the domain.

When comparing the machine translation with the human post-edit, BLEU scores for each language are between 84 and 90, and reference-based COMET scores for each language were between 92 and 94 (See Table 1 below). Of the 400 test translations reviewed and corrected by our reviewers, only 197 translations were modified (approximately 49%). These results present a pretty good baseline to start from, and will allow us, not only

to test the ability of LLMs at post-editing segments but also at leaving untouched those that do not need a correction.

Post-edition accuracy metrics - Baseline scores	
metric	raw_MT
DE	
BLEU	86.6
WER	11.87
COMET-REF	92.71
IT	
BLEU	84.6
WER	11.23
COMET-REF	93.59
PT	
BLEU	89.6
WER	9.16
COMET-REF	93.40
JA	
BLEU	86.1
WER	12.14
COMET-REF	93.51

Table 1: PE metrics of raw MT and human post-edited version

5.2 Results of the post-edition task

Table 3 presents the results of the post-edition task, categorized by metric and language.

Both Gemini and GPT demonstrated strong performance across all metrics, with Gemini outperforming in string-based automated evaluations (i.e., BLEU, MAD and WER). However, it is noteworthy that despite these promising performances, raw_MT consistently attained high scores across our chosen metrics, achieving the highest BLEU and WER scores across all languages and second best COMET-scores for German and Italian. This suggests that the introduced LLM post-editing did not significantly enhance the translation quality. A key factor contributing to this phenomenon might be the inclination of the models to excessively edit the machine-translated segment, resulting in deviations from the human post-editing.

Our experiments were conducted with the human post-edited version as the reference in metrics such as BLEU, WER, and MAD. Hence, it is plausible that while the LLM post-edit may not be inherently deficient, it may have overly altered the translation compared to the reference.

When considering the COMET QE score, which is calculated in a quality-estimation setting (i.e., without using a reference translation), we notice

that COMET assigns higher scores to LLM post-edits than the raw_MT and even the human post-edit. Further research should be carried out in order to understand whether these results are proof of an actual better quality from LLM outputs than the original raw MT, or whether COMET-QE might be biased towards machine-generated content.

Moreover, across languages, the raw machine translation also achieved the highest scores with string-based metrics, and outperformed most LLM responses when using reference-based COMET with the human post-edit as the reference. This result is consistent with our findings that LLMs might tend to make more extensive edits, leading to an increased divergence from the human post-edit.

While the models demonstrate capabilities in certain aspects, such as string-based assessments, their overall impact on translation quality requires further investigation. These findings emphasise the importance of refining post-editing strategies to align more effectively with human preferences and expectations.

5.3 Results of the MQM-analysis task

We now move to analyze the results obtained in the MQM task. In Table 4, we present all the metrics that were described in Section 4.3 including error detection accuracy, precision and recall (*ErrorAcc*, *ErrorPrecision* and *ErrorRecall*), error-type categorization accuracy (*ErrorTypeAcc*) and error-span detection precision (*ErrorSpanPrecision*).

In this regard, GPT seems to be the winner outperforming Claude and Gemini in all metrics provided, although not by a large margin. However, it is fair noticing that the *ErrorTypeAcc* (36.68%) and *ErrorSpanPrecision* are still quite low even for this model meaning that GPT shows promise at detecting sentences with errors but is still lagging behind at categorizing them according to the MQM types and detecting the actual error span.

MQM metrics			
metric	Claude	GPT_mqm	Gemini_mqm
ErrorAcc	65.75	66.75	64.5
ErrorPrecision	60.03	62.31	60.80
ErrorRecall	82.09	85.58	82.58
ErrorTypeAcc	35.17	36.68	36.68
ErrorSpanPrecision	22.07	18.23	9.79

Table 4: MQM accuracy metrics.

Error distribution of test data-set				
	Critical	Major	Minor	Total
Acc-Mistranslation	32%	44%	22%	92
Style	0%	7%	92%	68
Grammar	0%	24%	75%	45
Acc-Untranslated	10%	13%	75%	29
Domain	7%	46%	46%	28
Acc-Omission	0%	46%	53%	13
Typography	0%	9%	90%	11
Source	0%	14%	85%	7
Register	0%	20%	80%	5
Inconsistency	0%	25%	75%	4
Locale Convention	0%	0%	100%	4
Termbase	0%	0%	100%	4
Spelling	0%	33%	66%	3
Acc-Addition	0%	100%	0%	2
Unintelligible	0%	50%	50%	2

Table 2: Error distribution on the test data-set (sentences for post-edition + sentences for examples)

Post-edition accuracy metrics					
metric	human	raw_MT	claude_pe	GPT_pe	gemini_pe
DE					
BLEU	N/A	86.6*	74.42	70.65	77.55
MAD	N/A	N/A	11.30	<u>10.90</u>	8.99
WER	N/A	11.87*	22.73	24.27	<u>18.16</u>
COMET-REF	N/A	<u>92.71</u>	90.82	<u>91.71</u>	91.21
COMET-QE	42.38	<u>43.09</u>	43.42	42.8	42.92
IT					
BLEU	N/A	84.6*	72.79	<u>79.31</u>	65.92
MAD	N/A	N/A	<u>7.72</u>	7.97	6.58
WER	N/A	11.23*	22.81	21.93	<u>18.21</u>
COMET-REF	N/A	93.59	92.01	93.00	92.93
COMET-QE	36.43	<u>37.52</u>	38.49	35.59	37.33
PT					
BLEU	N/A	89.6*	79.86	<u>83.20</u>	75.60
MAD	N/A	N/A	<u>7.57</u>	9.01	5.79
WER	N/A	9.16	17.75	19.48	<u>13.32</u>
COMET-REF	N/A	93.40	92.56	92.72	<u>93.18</u>
COMET-QE	35.39	37.14	<u>38.36</u>	38.98	36.82
JA					
BLEU	N/A	86.1*	70.00	71.30	<u>76.13</u>
MAD	N/A	N/A	<u>11.67</u>	14.56	9.60
WER	N/A	12.14	19.46	22.02	14.91
COMET-REF	N/A	93.51	92.63	92.69	<u>93.45</u>
COMET-QE	<u>31.65</u>	31.02	31.6	33.52*	31.27

Table 3: Metrics for the PE methods and raw MT, with reference to the human post-edit. * indicates scores with a statistically significant difference from the second best score ($p < 0.05$).

For further exploring agreement between human annotations and models predictions, we calculate Cohen’s Kappa. Results in Table 5 show a similar behaviour to the aforementioned metrics: GPT_mqm shows the highest agreement with human annotators for MQM with a coefficient of 37.74. This is a compelling finding, since as, Popovic et. al (2014a) claimed in their research about inter-annotators’ agreement (IAA) on error-analysis, human annotators’ meta-understanding of language is variable, even when working with professional translators. In this paper the authors calculated IAA using Cohen’s Kappa in several languages. Their resulting coefficients were around 30 points for all the languages they studied.

Cohen’s Kappa Coefficient		
Claude_mqm	GPT_mqm	Gemini_mqm
35.76	37.74	30.17

Table 5: Cohen’s Kappa between human error annotations and predictions from models with the MQM prompt.

Finally, when exploring the resulting post-edited segments in this MQM setting, we found out that these tend to outperform those achieved by the PE prompt in many occasions (refer to Table 4 in Appendix B for a complete description of the metrics). However, while COMET-REF scores are higher for the MQM methods, string-based metrics are still higher for raw MT. In a similar manner as in the PE task, this suggests that LLMs are over-editing correct segments.

5.4 Results on the accuracy of models at selecting segments for post-edition

In Table 10 in Appendix C we offer a description of the number of segments which, according to linguists, needed a correction and those that were indeed corrected by the models (True Positives). This only includes the segments from the test-set (400 segments, 100 per language), since as it was mentioned in section 3.1, the extra 1,200 were passed as examples to the prompts. The most striking result here is that MQM methods correct 40% less than the PE methods, thus leading to a higher recall of the latter models. Depending on the production setting, this might be a desirable outcome where human review can be limited to the segments that were modified by the model. In a setting where a balance between precision and

recall is desired, GPT_pe was the best performing model with a f1-score of 70.66 points.

Table 6 highlights each LLM’s effectiveness in modifying segments containing errors as well as their ability to accurately modify the identified errors within those segments. To calculate the former, we just search for how many segments have been post-edited by the LLM. To calculate the latter, which is possibly more interesting for our research, we get the error span marked by the linguist and search for an exact match in the post-edited version. If the sub-string is not found in the post-edited sentence, we assume the error was modified. As shown in Table 6, we observe that Claude tends to modify more segments than the other two LLMs, and that the percentage of errors that were modified is below the percentage of modified segments. This points out once again that, while models are rewriting many segments, they are not always correcting the actual error that was marked by the linguist.

5.5 Qualitative analysis of generated MQM and post-editions

We further carry out a small manual analysis of the outputs of the LLMs in the quest for getting a better understanding of their behaviour. We decide to select Portuguese segments for its simplicity in analysis. Examining some of the responses from the PE and the MQM prompts, we observe the following:

- Sometimes the LLM detects the error and even gets the right type. However, while the PE prompt gets the post-edit right, the resulting fixed translation from the LLM using the MQM prompt is different from the one provided by the linguist (see Table 7). In the context of Life Sciences, there is a myriad of regulatory instructions as to how certain phrases should be translated, and while the LLM produces a correct translation it does not comply with the guidelines for this kind of documents. Adding a style guide in the prompt could the LLM produce a corrected version that follows the style and wordings from the guide.

Percentage of segments and errors modified by the LLMs per severity						
	Segments modified			Errors modified		
	Critical	Major	Minor	Critical	Major	Minor
Claude_mqm	80%	66%	84%	77%	45%	63%
Claude_pe	100%	100%	88%	91%	80%	56%
GPT_mqm	75%	61%	78%	66%	35%	64%
GPT_pe	94%	93%	95%	89%	85%	67%
Gemini_mqm	57%	43%	52%	54%	25%	38%
Gemini_pe	77%	81%	85%	69%	53%	49%

Table 6: Comparison of the LLMs' performance in modifying segments and errors per severity

Source: How does Nonafact work?
MT: Como funciona o não-afeto?
Human PE: Como o Nonafact funciona?
MQM:
Category: Accuracy - Mistranslation
Description: No relation to the source
Severity: Critical
Error span: não-afeto
GPT PE (PE prompt): Como o Nonafact funciona?
GPT PE (MQM prompt): Como funciona o Nonafact?
GPT MQM:
Category: Accuracy - Mistranslation
Severity: Critical
Description: -
Error span: não-afeto

Table 7: Example 1: Comparing human PE and MQM analysis with GPT's generated output

- In other occasions, we also observe that the LLM does not detect any error but still corrects the sentence, although in different ways depending on the prompt (see Table 8). Nevertheless, once again the post-editions do not align with the one provided by the linguist and seem to be changing parts of the sentence that were correct.

All these examples prove that there is still room for improvement on the use of LLMs for the tasks of post-edition and MQM analysis, specifically on the domain that has been researched throughout this paper. Firstly, by including a style guide of the preferred output and secondly, by encouraging a better alignment of the MQM analysis and the fixed translation generated.

Source: If the control of epilepsy is not clinically significantly improved after an adequate trial, vigabatrin treatment should not be continued.
MT: Se o controle da epilepsia não melhorar significativamente após um estudo clínico adequado, o tratamento com vigabatrina não deve ser continuado.
Human PE: Se o controle da epilepsia não melhorar significativamente após um estudo clínico adequado, o tratamento com vigabatrina não deverá ser continuado.
MQM:
Category: Fluency - Grammar
Description: Parallelism
Severity: Major
Error span: deve
GPT PE (PE prompt): Se o controle da epilepsia não apresentar uma melhoria clínica significativa após uma tentativa adequada, o tratamento com vigabatrina não deve ser continuado.
GPT PE (MQM prompt): Se o controle da epilepsia não melhorar significativamente após um período de teste adequado, o tratamento com vigabatrina não deve ser continuado.
GPT MQM: No errors found

Table 8: Example 2: Comparing human PE and MQM analysis with GPT's generated output

6 Discussion

After having carried out the automatic evaluation of the results obtained on the two proposed tasks (namely, APE and MQM analysis) we can conclude the following:

- With regards to APE, while there is still promise in using LLMs for improving MT outputs, as the COMET-QE scores from Table 3 suggest, when taking into account the compliance with a given reference segment,

LLMs do not seem to be “quite there” yet as other authors have previously pointed out (Kocmi et al., 2023). In order to find out whether LLMs did indeed improve the MT translations and that results from COMET-QE are not biased towards machine-generated outputs, further research should be carried out, for instance by obtaining pair-wise human preferences between translations.

- With regard to automatic MQM detection, while error-detection metrics present some promise, error-type categorization results are still only at 36% accuracy. While this metric is quite low, it is on par or slightly above reported inter-rater agreement in human evaluations such as that carried out by Popovic et al (2014a).
- In addition to error analysis, our MQM prompts also asked the models to produce a fixed translation. Comparing the fixed translations obtained in this way with those from the MTPE prompt, we observed that results generally improved for all languages and metrics, suggesting that the model benefits from the additional information and chain-of-thought style prompting. Further research could be carried out as to how removing the corrected translation from the examples given to the MQM prompt would affect these results.
- When comparing the accuracy of models at selecting segments for post-edition, we saw a large difference in the number of post-edited segments using the PE prompt vs. using the MQM prompt. The former tended to correct almost twice as much as the latter.
- Although, as we have mentioned, these models do not seem to be ready for production just yet, if there is an interest in using these models in completely independent workflow to carry out MQM analysis and PE, the choice as to which model to use should be made taking into account not only the accuracy of the edited content but the precision and recall metrics at selecting which segments indeed need to be post-edited as well, in order to reduce efforts while ensuring good results.
- Finally, when considering the results broken down by language, in general, we do not see

great variance across languages for any of the tasks. While the reference-based metrics for Japanese are often lower than for other languages, this is a common occurrence for this language. The commensurate performance across languages suggests that data contamination has not overly biased the results, and that the LLMs have strong priors for each of the languages we studied.

7 Future Work

Among our future work plans we intend to explore the fine-tuning of LLMs for the task of post-edition and MQM and compare the performance and costs with the approach proposed on this paper. Fine-tuning an LLM for a certain task has been proven to be a successful technique for achieving better results in certain tasks while reducing costs due to the shorter prompts that need to be sent to the model.

Another item of interest would be studying the integration of a style guide either by introducing it into the prompt or during fine-tuning. This could be useful for correction of stylistic errors for client customization.

Moreover, taking into account that our baseline models already achieve state-of-the-art performance, it would be interesting to carry out the same experiments on MT output which is objectively of poorer quality and analyze whether LLM post-edition and MQM analysis could significantly improve the translation in those cases.

Finally, human evaluation of the quality of LLM-post-edited content could be performed in order to get a better understanding of the results that were achieved with the automatic metrics presented on this paper.

References

- [Bechara et al.2011] Bechara, Hanna, Yanjun Ma, and Josef van Genabith. 2011. Statistical post-editing for a statistical MT system. In *Proceedings of Machine Translation Summit XIII: Papers*, Xiamen, China, September 19-23.
- [Bhattacharyya et al.2023] Bhattacharyya, Pushpak, Rajen Chatterjee, Markus Freitag, Diptesh Kanojia, Matteo Negri, and Marco Turchi. 2023. Findings of the WMT 2023 shared task on automatic post-editing. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*,

- pages 672–681, Singapore, December. Association for Computational Linguistics.
- [Cady et al.2023] Cady, Larry, Benjamin Tsou, and John Lee. 2023. Comparing Chinese-English MT performance involving ChatGPT and MT providers and the efficacy of AI mediated post-editing. In Yamada, Masaru and Felix do Carmo, editors, *Proceedings of Machine Translation Summit XIX, Vol. 2: Users Track*, pages 205–216, Macau SAR, China, September. Asia-Pacific Association for Machine Translation.
- [Correia and Martins2019] Correia, Gonalo M. and Andr  F. T. Martins. 2019. A simple and effective approach to automatic post-editing with transfer learning. In Korhonen, Anna, David Traum, and Llu s M rquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3050–3056, Florence, Italy, July. Association for Computational Linguistics.
- [do Carmo1 et al.2021] do Carmo1, F lix, Dimitar Shterionov, Joss Moorkens, Joachim Wagner, Murhaf Hossari, Eric Paquin, Dag Schmidtke, Declan Groves, and Andy Way. 2021. A review of the state-of-the-art in automatic post-editing.
- [Dyer et al.2013] Dyer, Chris, Victor Chahuneau, and Noah A. Smith. 2013. A simple, fast, and effective reparameterization of ibm model 2. In *In Proc. NAACL*.
- [Fernandes et al.2023] Fernandes, Patrick, Daniel Deutsch, Mara Finkelstein, Parker Riley, Andr  F. T. Martins, Graham Neubig, Ankush Garg, Jonathan H. Clark, Markus Freitag, and Orhan Firat. 2023. The devil is in the errors: Leveraging large language models for fine-grained machine translation evaluation.
- [Freitag et al.2021a] Freitag, Markus, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021a. Experts, Errors, and Context: A Large-Scale Study of Human Evaluation for Machine Translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474, 12.
- [Freitag et al.2021b] Freitag, Markus, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021b. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.
- [Gao et al.2023] Gao, Yuan, Ruili Wang, and Feng Hou. 2023. How to design translation prompts for chatgpt: An empirical study.
- [Garcia et al.2023] Garcia, Xavier, Yamini Bansal, Colin Cherry, George Foster, Maxim Krikun, Fangxiaoyu Feng, Melvin Johnson, and Orhan Firat. 2023. The unreasonable effectiveness of few-shot learning for machine translation.
- [Hendy et al.2023] Hendy, Amr, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation.
- [Junczys-Dowmunt and Grundkiewicz2016] Junczys-Dowmunt, Marcin and Roman Grundkiewicz. 2016. Log-linear combinations of monolingual and bilingual neural machine translation models for automatic post-editing. In Bojar, Ondr j, Christian Buck, Rajen Chatterjee, Christian Federmann, Liane Guillou, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Aur lie N v ol, Mariana Neves, Pavel Pecina, Martin Popel, Philipp Koehn, Christof Monz, Matteo Negri, Matt Post, Lucia Specia, Karin Verspoor, J rg Tiedemann, and Marco Turchi, editors, *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 751–758, Berlin, Germany, August. Association for Computational Linguistics.
- [Junczys-Dowmunt et al.2018] Junczys-Dowmunt, Marcin, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, Andr  F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia, July. Association for Computational Linguistics.
- [Kocmi and Federmann2023] Kocmi, Tom and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. In Nurminen, Mary, Judith Brenner, Maarit Koponen, Sirkku Latomaa, Mikhail Mikhailov, Frederike Schierl, Tharindu Ranasinghe, Eva Vanmassenhove, Sergi Alvarez Vidal, Nora Aranberri, Mara Nunziatini, Carla Parra Escart n, Mikel Forcada, Maja Popovic, Carolina Scarton, and Helena Moniz, editors, *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland, June. European Association for Machine Translation.
- [Kocmi et al.2023] Kocmi, Tom, Eleftherios Avramidis, Rachel Bawden, Ondr j Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamm Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popovi , and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore, December. Association for Computational Linguistics.

- [Lommel et al.2014a] Lommel, Arle, Maja Popovic, and Aljoscha Burchardt. 2014a. Assessing inter-annotator agreement for translation error annotation. In *MTE: Workshop on Automatic and Manual Metrics for Operational Translation Evaluation*, pages 31–37. Language Resources and Evaluation Conference Reykjavik.
- [Lommel et al.2014b] Lommel, Arle, Hans Uszkoreit, and Aljoscha Burchardt. 2014b. Multidimensional quality metrics (MQM): A framework for declaring and describing translation quality metrics. In *Tradumatica*, pages 455–463.
- [Lu et al.2024] Lu, Qingyu, Baopu Qiu, Liang Ding, Kanjian Zhang, Tom Kocmi, and Dacheng Tao. 2024. Error analysis prompting enables human-like translation evaluation in large language models.
- [McCann2020] McCann, Paul. 2020. fugashi, a tool for tokenizing japanese in python. *arXiv preprint arXiv:2010.06858*.
- [Moslem et al.2023] Moslem, Yasmin, Rejwanul Haque, John D. Kelleher, and Andy Way. 2023. Adaptive machine translation with large language models.
- [Pal et al.2017] Pal, Santanu, Sudip Kumar Naskar, Mihaela Vela, Qun Liu, and Josef van Genabith. 2017. Neural automatic post-editing using prior alignment and reranking. In Lapata, Mirella, Phil Blunsom, and Alexander Koller, editors, *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 349–355, Valencia, Spain, April. Association for Computational Linguistics.
- [Post2018] Post, Matt. 2018. A call for clarity in reporting bleu scores. *arXiv preprint arXiv:1804.08771*.
- [Raunak et al.2023] Raunak, Vikas, Amr Sharaf, Yiren Wang, Hany Hassan Awadallah, and Arul Menezes. 2023. Leveraging gpt-4 for automatic translation post-editing.
- [Rei et al.2022] Rei, Ricardo, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585.
- [Simard et al.2007] Simard, Michel, Cyril Goutte, and Pierre Isabelle. 2007. Statistical phrase-based post-editing. In Sidner, Candace, Tanja Schultz, Matthew Stone, and ChengXiang Zhai, editors, *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*, pages 508–515, Rochester, New York, April. Association for Computational Linguistics.
- [Tebbifakhr et al.2018] Tebbifakhr, Amirhossein, Ruchit Agrawal, Matteo Negri, and Marco Turchi. 2018. Multi-source transformer with combined losses for automatic post editing. In Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névéol, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 846–852, Belgium, Brussels, October. Association for Computational Linguistics.
- [Tiedemann2009] Tiedemann, Jörg, 2009. *News from OPUS - A Collection of Multilingual Parallel Corpora with Tools and Interfaces*, volume V, pages 237–248.
- [Vidal et al.2022] Vidal, Blanca, Albert Llorens, and Juan Alonso. 2022. In Campbell, Janice, Stephen Larocca, Jay Marciano, Konstantin Savenkov, and Alex Yanishevsky, editors, *Proceedings of the 15th Biennial Conference of the Association for Machine Translation in the Americas (Volume 2: Users and Providers Track and Government Track)*, pages 84–106, Orlando, USA, September. Association for Machine Translation in the Americas.
- [Vilar et al.2023] Vilar, David, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting palm for translation: Assessing strategies and performance.
- [Yang et al.2023] Yang, Xinyi, Runzhe Zhan, Derek F. Wong, Junchao Wu, and Lidia S. Chao. 2023. Human-in-the-loop machine translation with large language model.

8 Appendix A - Prompt Examples

8.1 GPT-MQM

”””Hello! Could you please tell me whether there is an error in the translation below?. Please, think it carefully before giving an answer:

If there is indeed an error or more than one error, could you please categorize them according to the error categories and subcategories below?

In order to carry out a proper analysis, first think of the error category and once you have that clear, subcategorize the error using the subcategories from that error category.

<error_category>Fluency: errors related to the linguistic well-formedness of the text, including problems with grammaticality, spelling, punctuation, and mechanical correctness.<error_category>

<sub_category>Domain Terminology: although the translation might be correct, it is not suited for the type of text.</sub_category>

<sub_category>Grammar: errors related to the grammar of a language</sub_category>

<sub_category>Inconsistency: translation is not consistent with previous words</sub_category>

<sub_category>Register: not using the right register (formal, informal, neutral)</sub_category>

<sub_category>Spelling: a term was misspelt, e.g., it contained to ss instead of one or a different letter was used.</sub_category>

<sub_category>Typography: typographic errors, related to punctuation or tags</sub_category>

<sub_category>Unintelligible: it is a made-up word or difficult to understand in normal language</sub_category>

<error_category>Terminology: errors arising when a term does not conform to normative domain or organizational terminology standards or when a term in the target text is not the correct, normative equivalent of the corresponding term in the source text.</error_category>

<sub_category>Domain Terminology: the error is a terminology issue deriving from the domain, the type of text: medical, tourism, daily life, law...</sub_category>

<error_category>Accuracy: errors occurring when the target text does not accurately correspond to the propositional content of the source text, introduced by distorting, omitting, or adding to the message.</error_category>

<sub_category>Addition: Addition of content.</sub_category>

<sub_category>Mistranslation: when a word has been translated differently that it should</sub_category>

<sub_category>Omission: omission of content</sub_category>

<sub_category>Untranslated: term was not translated</sub_category>

<error_category>Style: errors occurring in a text that are grammatically acceptable but are inappropriate because they deviate from organizational style guides or exhibit inappropriate language style.</error_category>

<error_category>Locale convention: errors

occurring when the translation product violates locale-specific content or formatting requirements for data elements.</error_category>

<error_category>Design: Errors regarding handling xml tags.</error_category>

<error_category>Source: There is an error on the SOURCE segment</error_category>.

Here are some examples that you can use as a reference:

Translation pair: {example}

Translation pair: {example}

Translation pair: {example}

Translation pair:
{translation pair to analyze and post-edit}
Analysis:””””

8.2 Gemini-PE

”””” As an expert linguist, your task is to perform post-editing (Light post-edit or Full post-edit) on machine-translated segments.

You will be working with {source language} as the source language and {target language} as the target language.

Below are three examples with human post-edits on the translations:

{example with source segment, translation, and post-edit}

{example with source segment, translation, and post-edit}

{example with source segment, translation, and post-edit}

Your task is to complete the following example by post-editing the translation, applying gender bias reduction if necessary. If no post-edit is needed, the post-edited translation should remain the same as the translation.

Example:
{example with source segment and translation}
””””

9 Appendix B - Post-edition results using the MQM prompt

Post-edition accuracy metrics (MQM prompt)					
metric	human	raw_MT	claude_mqm	GPT_mqm	gemini_mqm
DE					
BLEU	N/A	86.6	85.0	<u>85.1</u>	82.8
MAD	N/A	N/A	6.05	6.79	10.21
WER	N/A	11.87	<u>13.30</u>	13.56	37.40
COMET-REF	N/A	92.71	93.24	<u>93.10</u>	91.21
COMET-QE	42.38	43.09	43.61	40.33	42.92
IT					
BLEU	N/A	84.6*	79.90	82.90	<u>83.10</u>
MAD	N/A	N/A	4.99	<u>5.62</u>	6.20
WER	N/A	11.23	14.20	<u>12.91</u>	<u>12.91</u>
COMET-REF	N/A	<u>93.59</u>	79.55	93.95	79.71
COMET-QE	36.43	37.52	38.14	37.58	<u>38.15</u>
PT					
BLEU	N/A	89.6	88.80	89.50	88.50
MAD	N/A	N/A	4.46	5.65	<u>5.64</u>
WER	N/A	9.16	10.51	<u>9.95</u>	11.13
COMET-REF	N/A	93.40	93.86	93.95*	<u>93.42</u>
COMET-QE	35.39	37.14	37.78	36.91	40.14*
JA					
BLEU	N/A	86.1	82.90	<u>85.80</u>	81.70
MAD	N/A	N/A	<u>10.22</u>	8.62	11.82
WER	N/A	12.14	<u>14.88</u>	12.56	26.80
COMET-REF	N/A	93.51	93.20	93.55	91.83
COMET-QE	31.65	31.02	<u>33.2</u>	29.12	31.49

Table 9: Metrics for each of the methods and raw MT, with reference to the human post-edit. * indicates scores with a statistically significant difference from the second best score ($p < 0.05$)

10 Appendix C - Accuracy of models at choosing segments for post-edition

	Claude_mqm	Claude_pe	GPT_mqm	GPT_pe	Gemini_mqm	Gemini_pe
Needed Correction	197/400					
TP+FP	193	323	155	304	159	267
TP	131	183	111	177	86	161
percentage TP	66.50	92.89	56.35	89.85	43.65	81.73
percentage TN	69.46	31.03	78.33	37.44	64.04	47.78
percentage FN	33.50	7.11	43.65	10.15	56.35	18.27
precision	67.88	56.66	71.61	58.22	54.09	60.30
recall	66.50	92.89	56.35	89.85	43.65	81.73
f1-score	67.18	70.38	63.07	70.66	48.31	69.40

Table 10: Accuracy of models at choosing which segments to post-edit. If a segment needed a correction and was post-edited it is counted as a True Positive, while if a segment did not need a correction and was left untouched, it is counted as a True Negative

Post-editors as Gatekeepers of Lexical and Syntactic Diversity: Comparative Analysis of Human Translation and Post-editing in Professional Settings

Lise Volkart

FTI/TIM, University of Geneva
Switzerland
lise.volkart@unige.ch

Pierrette Bouillon

FTI/TIM, University of Geneva
Switzerland
pierrette.bouillon@unige.ch

Abstract

This paper presents a comparative analysis between human translation (HT) and post-edited machine translation (PEMT) from a lexical and syntactic perspective to verify whether the tendency of neural machine translation (NMT) systems to produce lexically and syntactically poorer translations shines through after post-editing (PE). The analysis focuses on three datasets collected in professional contexts containing translations from English into French and German into French. Through a comparison of word translation entropy (HTRa) scores, we observe a lower degree of lexical diversity in PEMT compared to HT. Additionally, metrics of syntactic equivalence indicate that PEMT is more likely to mirror the syntactic structure of the source text in contrast to HT. By incorporating raw machine translation (MT) output into our analysis, we underline the important role post-editors play in adding lexical and syntactic diversity to MT output. Our findings provide relevant input for MT users and decision-makers in language services as well as for MT and PE trainers and advisers.

1 Introduction

Post-editing (PE) has now largely proved to be a good alternative to purely human translation (HT) in professional contexts. By allowing certain productivity gains without negatively affecting

the quality of the final translation (Daems, 2016; Läubli et al., 2019), MT and PE have found their place in professional translation workflows. Nevertheless, translators often express mixed feelings towards MT. On one hand, the tool is appreciated for its help when dealing with high workloads and time constraints, but on the other hand, it is perceived as a threat to translation's creativity, originality and naturalness (Alvarez-Vidal et al., 2020; Girletti, 2024). These concerns are legitimate: numerous studies have revealed the NMT tendency to produce an output that is less lexically varied and syntactically closer to source text than HT (Vanmassenhove et al., 2019; Toral, 2019; Vanmassenhove et al., 2021; Webster et al., 2020; Ahrenberg, 2017; Shaitarova et al., 2023; Luo et al., 2024). Furthermore, some studies have found measurable differences on parallel corpora of HT and post-edited machine translation (PEMT) in terms of lexical diversity, lexical density and sentence length, among others, suggesting the existence of a *post-editese* phenomenon (Castilho et al., 2019; Castilho and Resende, 2022; Toral, 2019). However, Volkart and Bouillon (2023), demonstrated the difficulty of generalising these findings over different corpora, domains and language pairs, particularly when analysing authentic comparable corpora. Such corpora, in which HT and PEMT are the translations of different source texts, necessitate analysis with metrics that encompass the attributes of the source. Although demanding, the study of authentic HT and PEMT corpora is crucial for developing a detailed understanding of the distinct characteristics of PEMT output in professional contexts.

In this study, we compare the lexical and syntactic characteristics of authentic HT and PEMT output produced in professional contexts relying

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

on metrics suitable for the study of comparable corpora. Whenever possible, raw MT output is added to the analysis to confirm initial premises on raw MT, as well as findings by previous studies. We measure the variety of translation solutions by automatically computing average word translation entropy (HTra) (Carl et al., 2016) and the syntactic equivalence between source and target based on three syntactic equivalence metrics from the AS-TRed library (Vanroy et al., 2021) to answer the following research question:

Is the NMT tendency toward lexical impoverishment and source sentence structure mirroring discernible in the PEMT final product?

Throughout our analysis, we make the following observations:

1. PEMT final output is affected by NMT bias in terms of lexical diversity and presents lower levels of translation variety
2. Syntactic shining through from raw NMT occurs in PEMT output but remains limited
3. PE adds significant levels of lexical and syntactic variety to MT output

By including three authentic corpora, two language pairs, various state-of-the-art NMT systems and carefully selected metrics, our study contributes to improve our understanding of the impact of MT integration on translated language. Our findings provide valuable insights to inform decisions about where and when to use or not to use MT and can contribute to enhancing PE training programs and refining best practices. Finally, the comparison between HT, PEMT and raw MT reaffirms the essential role played by human post-editors.

Section 2 gives a brief overview of previous relevant studies. Datasets and corpora are described in Section 3 and experimental setup in Section 4. We present and comment our results in Section 5. Section 6 briefly presents illustrative examples. Finally, Section 7 reports our conclusions and main findings.

2 Related work

Lexical level Vanmassenhove et al. (2019), compared statistical machine translation (SMT), NMT and HT in terms of lexical richness, to verify the hypothesis according to which data-driven MT

systems, due to their probabilistic nature, would tend to favour more frequent words and disregard less frequent ones and therefore produce a lexically less diverse output than HT. Their experiment, conducted on 12 different MT systems (SMT and NMT with different architectures and with and without using backtranslated data) and two languages directions (EN-FR and EN-ES), confirmed this hypothesis. The three investigated lexical richness metrics (Yule's I, type/token ratio and measure of textual lexical diversity) indicated a lower lexical richness in the MT output in comparison to HT. With a further analysis on word frequencies, the authors demonstrated the tendency of MT systems to increase the frequency of already frequent words while decreasing the frequency of less frequent ones. This tendency toward overgeneralisation was again observed by Vanmassenhove et al. (2021) when measuring the difference in lexical and morphological richness between an MT system's training data and its output for the language directions EN-FR and EN-ES. They measured a loss of lexical and morphological richness between the training data and the system's output. Webster and al. (2020) also observed a loss of lexical richness and a homogenisation of lexicon with NMT when comparing lexical richness of literary excerpts translated by humans and by two online NMT systems. Although pursuing a slightly different goal (i.e. comparing raw MT, PEMT and revised PEMT), the work by Macken et al. (2022) is worth mentioning here, particularly because the authors relied, among others, on the average automatic word translation entropy (denoted AWTE in their paper) to assess lexical richness of the different translation modes. Overall, their experiment showed that PE and revision tend to increase the lexical variety of the raw MT, with AWTE being the most unequivocal of the three metrics used (AWTE, TTR and Mass Index).

Syntactic level As for the syntactic profile of MT outputs, several studies investigated the syntactic similarity between source and target for HT and PEMT. In 2017, Ahrenberg (2017) found out that, when translating from English into Swedish, NMT tends to produce an output that mimic the source structure, performing less word re-ordering than human translators. Comparing HT and generic NMT on literary excerpts with the help of word-cross and AS-TRed metrics (Vanroy et al., 2021), Webster et al. (2020) found out NMT tends to re-

main syntactically closer to the source structure. The same tendency was observed by Shaitarova et al. (2023), who tested syntactic equivalence between source and target using the ASTrED tool to compute cross-alignments on several large corpora. Their comparison of HT and NMT from different commercial systems indicated a general tendency of NMT systems to reproduce the syntax of the source, whereas HT appears to be more creative on this aspect. It is worth noting that in this experiment, out of the 4 tested NMT systems, DeepL appeared as the one producing the most syntactically diversified output. Finally, in an extensive study comparing NMT and HT in terms of morphosyntactic divergence between source and target on three language directions, Luo et al. (2024) found out that NMT tends to produce less diverse morphosyntactic patterns and more one-to-one alignments than HT.

3 Datasets

Our experiment is based on three authentic datasets containing professional translations collected from in-house language services¹. Each dataset contains a balanced amount of HT and PEMT segments with their respective source. Dataset ENfr1 was compiled from the same data as in Volkart and Bouillon (2022) and contains translations from English into French extracted from documents of the European Investment Bank (EIB). Dataset ENfr2 and dataset DEfr are derived from the dataset described in Volkart and Bouillon (2023). Dataset ENfr2 contains translations from English into French shared with us by a sports organisation based in Switzerland, while dataset DEfr contains translations from German into French collected from an insurance company. For all datasets, raw MT used for PE came from various state-of-the-art NMT systems (generic and/or customised). Table 1 presents the size of each dataset and corpus. In addition to PEMT data, we added raw MT of the PEMT source data to ENfr1 and ENfr2 datasets. This raw MT was generated for this experiment using DeepL Pro² ³. Original raw MT is not saved by the language services during the PE process, which

¹All services shared their data on a voluntary basis. Agreements between researchers and organisations were signed when needed and data was anonymized when required.

²in february 2024

³Data provider of the DEfr corpus did not allow us to translate their corpus using an online MT system

restricts our analysis of authentic data to PEMT product. However, we deem informative to include an example of raw MT output, although artificial, in our analysis. It allows us, among others, to verify if the tendencies observed by previous studies on the lexical and syntactic profile of MT outputs are indeed to be seen in our data.

4 Experiment

4.1 Variety of translation solutions

To measure the variety of translation solutions, we rely on the Word translation entropy metric (denoted HTra) (Carl et al., 2016). HTra is computed as the sum over all observed word translation probabilities $p(s \rightarrow t_i)$ of a given source text word s into target text word $t_i \dots n$ multiplied with their information content $I(p) = -\log_2(p)$ (Carl et al., 2016) as is the following equation:

$$HTra(s) = -\sum_{i=1}^n p(s \rightarrow t_i) \times \log_2(p(s \rightarrow t_i))$$

This score reflects, for a given source word, the amount of translation alternatives and their distribution in the target (Bangalore et al., 2016; Gilbert et al., 2023). The higher the HTra, the higher the variety in the translation of that source word in the target corpus. Compared to the TTR-based scores (such as TTR (Scott, 2019), STTR (Scott, 2019), MSTTR (Malvern and Richards, 2002) or MATTR (Covington and McFall, 2010)) often used to compare lexical richness of HT and MT/PEMT, HTra offers two main advantages: first, it is computed on the target corpus given its source and therefore allows us to compare translations from different source texts more easily (whereas TTR requires us to take into account the influence of the source while comparing the target, see Volkart and Bouillon (2022) and Volkart and Bouillon (2023) for a more detailed discussion on this aspect), and second, it encompasses two different aspects of lexical/translation richness that are the number of unique translation solutions, and also their distribution (does one solution account for 90% of the occurrences or are all translation solutions equally used by the translator?). Then, in addition to being more appropriate regarding our corpus design, HTra captures more information on the lexical richness of different translations than TTR-based metrics.

Computing HTra for a given source word requires the extraction all occurrences of that source

Dataset	Trans. mode	# segment pairs	# source tokens
ENfr1	HT	1,852	40,560
	PEMT	1,852	41,803
ENfr2	HT	2,280	43,379
	PEMT	2,280	49,896
DEfr1	HT	7,769	106,864
	PEMT	7,769	106,673

Table 1: Number of segments and source tokens for each dataset and translation mode

word with their respective translations in the target corpus. Whereas it can be done manually on small corpora or for a selection of source words such as in Volkart and Bouillon (2022), where it was computed for a set of 20 adverbs, it can rapidly become impossible to apply on large corpora.

We computed HTra automatically using ad hoc python scripts. Word alignment was performed with awesome-align (Dou and Neubig, 2021), a neural word aligner based on multilingual BERT, without fine-tuning. Out of this automatic alignment, we extracted the list of source-target pairs for content words (adverbs, adjectives, nouns and verbs) and grouped source words aligned with multiple target words together to form one-to-many alignments. Non-aligned source words were added to the list as non-translated. Tagging and lemmatization were performed in parallel using SpaCy’s transformer models for English, French and German ⁴. We computed HTra for all content source lemmas that occur at least three times in both HT and PEMT source corpora.

This automatic HTra computation pipeline was validated against manually computed scores from Volkart and Bouillon, on the same corpus and the same subset of adverbs (2022). Pearson’s correlation coefficients between automatic and manual HTra scores are respectively 0,83 for HT and 0,81 for PEMT. These high levels of correlation validate the automatic calculation method as well as the quality of the automatic word alignment.

4.2 Syntactic equivalence

To measure the impact of PEMT on the syntactic level, we used the ASTrED python library (Vanroy et al., 2021) to compute three metrics of syntactic equivalence, namely the label changes, the Syntactically Aware Cross (SACr) and the Aligned syntactic tree edit distance (ASTrED). Those metrics aim at capturing syntactic equivalence

between a source and a target segment based on differences in word/word group order, differences in dependency labels and differences in syntactic structures (Vanroy et al., 2021). The ASTrED library relies on Stanza parser (Qi et al., 2020) for universal dependency parsing and an adapted version of awesome-align (Dou and Neubig, 2021) for word alignment (Vanroy et al., 2021).

Label changes Label changes correspond, for a given source-target sentence pair, to the number of source-target word-aligned pairs that have different dependency labels, normalised by the total number of alignments for that sentence. This metric captures the linguistic differences between aligned words on the surface level (Vanroy et al., 2021).

SACr SACr quantifies the degree of reordering of word sequences that occurred between source and target (Vanroy et al., 2021). Words are grouped together according to their relation in the dependency tree to form linguistically motivated word sequences. Source and target word sequences are then aligned based on word alignments and SACr value is computed by dividing the number of cross-alignments normalised by the total number of alignments. SACr captures the surface word order differences between the source and target sentences.

ASTrED ASTrED captures the source and target structural differences on a deeper level by comparing dependency trees while taking word alignments into account. The computed tree edit distance is normalised by the average number of source and target words (Vanroy et al., 2021). For further details and illustrated examples on these metrics, we invite the reader to refer to Vanroy et al. (2021).

⁴<https://spacy.io/models>

5 Results

5.1 Lexical richness

HTra was automatically computed on content lemmas for HT, PEMT and raw MT⁵. Average scores by POS categories and for all content lemmas for HT, PEMT and raw MT corpora are presented in Table 2.

All POS categories together, the first thing we observe is that the raw MT generated with DeepL for this experiment presents a much lower variety of translation solutions compared to HT. For all categories together, the HTra score is more than 20% lower for raw MT. This confirms what has been observed in previous studies regarding the tendency of MT systems to narrow the range of translation solutions by increasing the frequency of already frequent words while decreasing the frequency of less frequent ones. This tendency shines through in the PEMT output which exhibits a generally lower HTra score compared to HT for all three datasets. Datasets ENfr2 and DEfr present very similar results, with HTra almost 8% lower for PEMT in contrast to HT. This loss of translation solution variety is less marked in the ENfr1 dataset, but still to be seen.

Those scores indicate that post-editors presumably add significant amount of lexical variety to the MT output, but still not enough to reach the level of variation from HT.

Looking at HTra scores by POS category separately, we see that the observed loss of translation variety is spread differently across categories for the different datasets. For ENfr1, adverbs show the biggest loss of variety in PEMT, while this loss is very limited for verbs. For ENfr2, on the contrary, the loss of variety affects primarily nouns and verbs, while for adverbs we even observe a higher translation variety in PEMT. Interestingly, the loss of translation solution variety in raw MT seems to correlate with the loss of translation solution variety for PEMT for this dataset. Finally, in DEfr, the loss of variety in PEMT is more evenly spread across categories, with a slightly stronger effect for adjectives and adverbs. Here, different POS categories appear to be differently affected by the loss of translation variety in PEMT depending on the dataset and, presumably, on the language

⁵To prevent the results from being overly biased by non-frequent or topic-related lemmas, HTra was computed for content lemmas occurring at least three times in both source corpora

pair.

5.2 Syntactic equivalence

Table 3 presents the average scores for HT for all three metrics and the relative differences for PEMT and raw MT when available. For all three metrics, a lower score indicates a higher level of syntactic equivalence between source and target. Similarly to what we observe for HTra, raw MT differs significantly from HT for both English into French datasets, with all three metrics indicating that the target tends to be syntactically closer to the source for raw MT. Once again, it is coherent with what could be observed in other studies. The difference between HT and PEMT is less straightforward. For the ENfr1 dataset, word sequence reordering, as measured by SACr, and label changes are more frequent in PEMT, whereas tree edit distance is slightly lower for PEMT. As for ENfr2, metrics show that PEMT is syntactically closer to source than HT, with significantly less label changes and lower tree edit distance. As for DEfr, SACr and label changes are not significantly different for PEMT and HT, but ASTrED points toward more similarity between source and target on the deeper level in PEMT. These results show that PE clearly blurs the line between HT and MT on the syntactic level. Post-editors play a major role in adding syntactic variety to the MT output during post-editing especially on the surface level by adding large amounts of word reordering and dependency label change. On a deeper level however, PEMT stays closer to the source than HT as expressed by consistently lower ASTrED scores.

6 Examples

Loss of translation solution variety: to illustrate what a lower HTra score concretely means, we present examples showing the translation solutions distribution for particular lemmas in Figures 1 and 2 in Appendix A. Figure 1 shows the distribution of translation solutions for the noun “impact” in HT, PEMT and raw MT within the ENfr2 dataset. While PEMT and HT exhibit an equal number of translation solutions, their frequencies are more evenly dispersed in HT, resulting in a HTra score of 2.84 for HT compared to 2.79 for PEMT. In contrast, raw MT yields only two distinct translation solutions (with the absence of translation considered as a translation “choice”), where one solution overwhelmingly dominates,

Corpus	ENfr1			ENfr2			DEfr	
	HT	PEMT	RawMT	HT	PEMT	RawMT	HT	PEMT
ADJ	1.32	-1.52%	-17.42%*	1.45	-2.76%	-20.00% [◊]	1.69	-9.47% [◊]
ADV	1.51	-4.64%	-15.89%	1.69	+2.37%	-10.65%	1.93	-9.33% [◊]
NOUN	1.07	-2.80%	-24.30% [◊]	1.26	-10.32% [◊]	-33.33% [◊]	1.16	-6.90% [◊]
VERB	1.92	-0.52%	-23.44% [◊]	2.16	-9.72% [◊]	-23.15% [◊]	1.90	-7.89% [◊]
All	1.34	-1.49%	-21.64%[◊]	1.54	-7.79%[◊]	-25.97%[◊]	1.53	-7.84%[◊]

Table 2: Average HTra scores for HT and relative difference for PEMT and raw MT for all content lemmas and each POS category. *indicate significance at $p < 0.005$ and [◊] at $p < 0.001$. Significance was tested using Mann Whitney non-parametric test.

Corpus	ENfr1			ENfr2			DEfr	
	HT	PEMT	RawMT	HT	PEMT	RawMT	HT	PEMT
ASTrED	0.6153	-1.67%	-3.79% [◊]	0.6348	-7.29% [◊]	-14.07% [◊]	0.6518	-3.01% [◊]
SACr	0.2702	+6.62% [◊]	-24.46%*	0.2558	-3.36%	-40.89% [◊]	0.3185	-0.78%
Label Ch.	0.1982	+1.11%	-5.90% [◊]	0.2186	-7.55% [◊]	-11.89% [◊]	0.2292	+0.26%

Table 3: Syntactic equivalence scores for HT and relative difference for PEMT and raw MT. *indicate significance at $p < 0.005$ and [◊] at $p < 0.001$. Significance was tested using Mann Whitney non-parametric test

resulting in an HTra score of 0.24. Looking at the adverb “also” within the ENfr1 dataset presented in Figure 2, we note that in this situation the high HTra score for HT (1.27) is principally due to the number of different translation solutions, more than to their frequency distribution. PEMT achieves an HTra score of only 1.02 while raw MT, due to the strong dominance of the most frequent solution barely reaches 0.85. These examples show how the loss of lexical diversity occurs in PEMT through the loss of translation solution variety and how the use of MT can lead to a reinforcement of the most frequent translation solutions at the expense of the less frequent ones.

Syntactic equivalence: Figures 3 and 4 show the word alignments between source and target for two sentences extracted from our datasets. They illustrate two contrasting examples regarding syntactic equivalence between source and target. In Figure 3, the target sentence presents a high level of syntactic equivalence according to the three computed scores (ASTrED = 0.29, SACr = 0.05, Label change = 0.16) and this is intuitively expressed in the word alignment. The target sentence is an almost one-to-one translation of the source with minimal word reordering and dependency label changes. Figures 4 in contrast presents the word alignments with higher scores and therefore less syntactic equivalence (ASTrED = 0.78, SACr = 0.58, Label change = 0.25). Here again, just by looking at the word alignment,

it is clear that the target presents higher levels of word reordering and structural differences compared to the source sentence.

7 Conclusion

This paper compares authentic sets of HT and PEMT produced in three different professional contexts for the language directions English into French and German into French, with additional analysis incorporating raw MT output from DeepL for two of the datasets. The objective is to compare HT and PEMT in terms of lexical and syntactic variety to verify whether the general tendency of NMT systems to produce lexically and syntactically less varied output still shines through after a PE step performed by professional translators. Using HTra for the lexical aspects and ASTRaED, SACr and dependency label changes (Vanroy et al., 2021) for the syntactic aspects, we note the strong tendency of raw NMT (in this case DeepL) to produce lexically less varied translations that tend to mirror the source sentence structure. This tendency is strongly attenuated by the PE step, with PEMT output being generally closer to HT than to raw MT on both aspects. This indicates that post-editors presumably add significant levels of lexical and syntactic variety to the MT output (“presumably”, because raw MT under analysis is not the one originally used for PEMT, but we assume it reflects the general level of lexical and syntactic variation of NMT systems). Still, the final PEMT output does not systematically reach the same level

of variety as HT, especially on the lexical level. For all datasets PEMT exhibits lower levels of translation solution variety. A tendency towards more syntactic equivalence between source and target in PEMT is clear for one dataset but more nuanced for the two others. These findings are particularly relevant in contexts where lexical and syntactic variety are regarded as criteria for assessing translation quality.

Furthermore, our work highlights the crucial importance of PE, not only in ensuring the accuracy of the target text, but also in maintaining an adequate level of lexical diversity and syntactic naturalness in the final translation. While this aspect may seem unimportant for certain types of texts, it holds significant relevance for others. In many cases, (human) translation is not only about overcoming language barriers but also about producing “texts that satisfy the linguistic norms of a target culture and are adapted to the assumed knowledge of its reader” (Ahrenberg, 2017, 1). It is also of utmost importance considering the fact that PEMT output is likely to be re-used to train NMT systems, and therefore to amplify over and over the already existing biases.

Finally, we emphasise the relevance of our findings for the improvement of post-editing training programs and guidelines. While translators are still today often advised to stick to the TAUS PE guidelines (TAUS and CNGL, 2010) and to not intervene on the stylistic level, we are convinced that adding lexical and syntactic diversity (even when not strictly necessary from micro-level perspective) to MT output is essential to preserve the quality of translated text at the macro-level.

Acknowledgement

We would like to thank the language services who kindly accepted to share their translation memories for this research project. We would also like to thank the reviewers for their insightful comments and feedback.

References

Ahrenberg, Lars. 2017. Comparing machine translation and human translation: A case study. In *RANLP 2017: The First Workshop on Human-Informed Translation and Interpreting Technology (HiT-IT)*, pages 21–28. Association for Computational Linguistics.

Alvarez-Vidal, Sergi, Antoni Oliver, and Toni Ba-

dia. 2020. Post-editing for professional translators: cheer or fear? *Tradumàtica*, (18):0049–69.

- Bangalore, Srinivas, Bergljot Behrens, Michael Carl, Maheshwar Ghankot, Arndt Heilmann, Jean Nitzke, Moritz Schaeffer, and Annegret Sturm. 2016. Syntactic variance and priming effects in translation. In *New directions in empirical translation process research*, pages 211–238. Springer.
- Carl, Michael, Moritz Schaeffer, and Srinivas Bangalore. 2016. The CRITT translation process research database. In *New directions in empirical translation process research*, pages 13–54. Springer.
- Castilho, Sheila and Natália Resende. 2022. Post-editeuse in literary translations. *Information*, 13(2):66. Publisher: Multidisciplinary Digital Publishing Institute.
- Castilho, Sheila, Natália Resende, and Ruslan Mitkov. 2019. What influences the features of post-editeuse? a preliminary study. In *Proceedings of the Human-Informed Translation and Interpreting Technology Workshop (HiT-IT 2019)*, pages 19–27. Varna, Bulgaria.
- Covington, Michael A. and Joe D. McFall. 2010. Cutting the Gordian Knot: The Moving-Average Type-Token Ratio (MATTR). *Journal of Quantitative Linguistics*, 17(2):94–100.
- Daems, Joke. 2016. *A translation robot for each translator?: A comparative study of manual translation and post-editing of machine translations: Process, quality and translator attitude*. PhD Thesis, Ghent University.
- Dou, Zi-Yi and Graham Neubig. 2021. Word alignment by fine-tuning embeddings on parallel corpora. In *Conference of the European Chapter of the Association for Computational Linguistics (EACL)*.
- Gilbert, Devin, Cristina Toledo-Báez, Michael Carl, and Haydeé Espino. 2023. Impact of word alignment on word translation entropy and other metrics. In Lacruz, Isabel, editor, *Translation in Transition: Human and machine intelligence*, page 203. Publisher: John Benjamins Publishing Company.
- Girletti, Sabrina. 2024. *Working with Pre-translated Texts: Investigating Machine Translation Post-editing and Human Translation Revision at Swiss Corporate In-house Language Services*. Ph.D. thesis, University of Geneva.
- Läubli, Samuel, Chantal Amrhein, Patrick Düggelin, Beatriz Gonzalez, Alena Zwahlen, and Martin Volk. 2019. Post-editing productivity with neural machine translation: An empirical assessment of speed and quality in the banking and finance domain. In *Proceedings of Machine Translation Summit XVII: Research Track*, pages 267–272, Dublin, Ireland. European Association for Machine Translation.

- Luo, Jiaming, Colin Cherry, and George Foster. 2024. To Diverge or Not to Diverge: A Morphosyntactic Perspective on Machine Translation vs Human Translation. *arXiv preprint arXiv:2401.01419*.
- Macken, Lieve, Bram Vanroy, Luca Desmet, and Arda Tezcan. 2022. Literary translation as a three-stage process: machine translation, post-editing and revision. In *23rd Annual Conference of the European Association for Machine Translation*, pages 101–110. European Association for Machine Translation.
- Malvern, David and Brian Richards. 2002. Investigating accommodation in language proficiency interviews using a new measure of lexical diversity. *Language Testing*, 19(1).
- Qi, Peng, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.
- Scott, Mike. 2019. WordSmith tools manual.
- Shaitarova, Anastassia, Anne Göhring, and Martin Volk. 2023. Machine vs. Human: Exploring Syntax and Lexicon in German Translations, with a Spotlight on Anglicisms. In *The 24rd Nordic Conference on Computational Linguistics*.
- TAUS and CNGL. 2010. Machine Translation Post-Editing Guidelines.
- Toral, Antonio. 2019. Post-editeese: an exacerbated translationese. In *Proceedings of MT Summit XVII*, volume 1, pages 273 – 281. Dublin, Ireland.
- Vanmassenhove, Eva, Dimitar Shterionov, and Andy Way. 2019. Lost in translation: Loss and decay of linguistic richness in machine translation. In *Proceedings of MT Summit XVII*, volume 1, pages 222 – 232. Dublin, Ireland.
- Vanmassenhove, Eva, Dimitar Shterionov, and Matthew Gwilliam. 2021. Machine Translationese: Effects of Algorithmic Bias on Linguistic Complexity in Machine Translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2203–2213, Online, April. Association for Computational Linguistics.
- Vanroy, Bram, Orphée De Clercq, Arda Tezcan, Joke Daems, and Lieve Macken. 2021. Metrics of syntactic equivalence to assess translation difficulty. In *Explorations in empirical translation process research*, pages 259–294. Springer.
- Volkart, Lise and Pierrette Bouillon. 2022. Studying Post-Editese in a Professional Context: A Pilot Study. In *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, pages 71–79.
- Volkart, Lise and Pierrette Bouillon. 2023. Are post-editeese features really universal? In Orăsan, Constantin, Ruslan Mitkov, Gloria Corpas Pastor, and Johanna Monti, editors, *Proceedings of the International Conference on Human-Informed Translation and Interpreting Technology (HiT-IT 2023)*, pages 294–304, Naples.
- Webster, Rebecca, Margot Fonteyne, Arda Tezcan, Lieve Macken, and Joke Daems. 2020. Gutenberg goes neural: Comparing features of dutch human translations with raw neural machine translation outputs in a corpus of english literary classics. In *Informatics*, volume 7, page 32. MDPI. Issue: 3.

Appendix A. Examples

HT		PEMT		Raw MT	
effet (8)		effet (6)		impact (25)	
impact (5)		impact (3)		not translated (1)	
répercussion (5)		répercussion (6)		(0)	
influence (4)		force (1)		(0)	
évaluation (3)		action (1)		(0)	
not translated (2)		incidence (1)		(0)	
valeur ajoutée (1)		conséquence (1)		(0)	
retombée (1)		retombée (1)		(0)	
évaluation (1)		différence (1)		(0)	

Figure 1: Translation solutions distribution for the noun “impact” in ENfr2 HT, PEMT and raw MT.

HT		PEMT		Raw MT	
également (76)		également (60)		également (67)	
aussi (30)		aussi (12)		aussi (7)	
outre (3)		que (1)		pour (1)	
parallèle (1)		ou (1)		not translated (1)	
que (1)		notamment (1)		tout (1)	
ailleurs (1)		y compris (1)		y compris (1)	
remercier (1)		(0)		(0)	
autre (1)		(0)		(0)	

Figure 2: Translation solution distribution for the adverb “also” in ENfr1 HT, PEMT and raw MT.

Analysing your own actions and the feedback you receive from other referees is useful for preparing for future matches and growing your knowledge base .

Analysier vos propres actions et les commentaires que vous recevez de la part d' autres arbitres est utile pour préparer vos futurs matches et développer votre base de connaissances .

Figure 3: Example of a source-target sentence pair presenting high levels of syntactic equivalence, with automatic word alignments.

We have seen incredible demand so far .

Jusqu' ici , la demande a été incroyable .

Figure 4: Example of a source-target sentence pair presenting low levels of syntactic equivalence, with automatic word alignments.

Exploring NMT Explainability for Translators Using NMT Visualising Tools

Gabriela Gonzalez-Saez¹, Mariam Nakhle^{1 5}, James Robert Turner⁴,
Fabien Lopez¹, Nicolas Ballier³, Marco Dinarelli¹, Emmanuelle Esperança-Rodier¹,
Sui He⁴, Raheel Qader⁵, Caroline Rossi², Didier Schwab¹, Jun Yang⁴

¹Univ. Grenoble Alpes, CNRS, Grenoble INP, LIG 38000 Grenoble, France; ²Université Grenoble Alpes; ³Université Paris Cité, LLF & CLILLAC-ARP, 75013 Paris, France;

⁴Swansea University; ⁵Lingua Custodia, France

`gabriela-nicole.gonzalez-saez@univ-grenoble-alpes.fr`

Abstract

This paper describes work in progress on Visualisation tools to foster collaborations between translators and computational scientists. We aim to describe how visualisation features can be used to explain translation and NMT outputs. We tested several visualisation functionalities with three NMT models based on Chinese-English, Spanish-English and French-English language pairs. We created three demos containing different visualisation tools and analysed them within the framework of performance-explainability, focusing on the translator's perspective.

1 Introduction

The development of machine translation (MT) is influenced by a wide range of actors and agents, ranging from the investors to general public. A stakeholder approach to MT enables us to examine the effects of MT on each of the different interest groups, with particular reference to levels of involvement with MT (e.g., translators, students and trainees, end users, MT investors and developers, translation agencies, and academic researchers) (Guerberof-Arenas and Moorkens, 2023).

Upon refining the landscape of MT to include the directly associated stakeholders, several distinct categories emerge. The primary category consists of MT developers, typically computer scientists, whose focus lies in enhancing the accuracy and fluency of translations. In contrast, a second group of stakeholders, comprised of linguists

and translators with expertise in translation studies, may argue that translation quality cannot be regarded as a static or absolute concept; instead, it is influenced by both subjective and objective factors that may change over time, as evidenced by semantic/communicative and functional translation approaches. Their concern centres on how MT fits into their practical translation workflow, and MT's ability to handle cultural-specific items and nuances – a realm in which translators take great pride. Moreover, industry surveys, such as the Freelance Translator Survey 2023 by Inbox Translation, and CIOL Insights 2022,¹ have also demonstrated that translators are primarily concerned about the effects of MT in their professional status, i.e., decreased translation/post-editing rates, clients' unrealistic expectations and other people's perception of their professionalism. Moreover, end-users constitute another critical group, prioritising usability, speed, and the cost-effectiveness of translations (Vieira et al., 2023).

Although these three groups of stakeholders are closely related to the development of MT, they do not always understand each other's work or demands, underscoring the need for continuous dialogue between each group. This work is of part of the MAKE-NMTViz project, which, by bringing together key stakeholders in MT development, aims to connect MT researchers with professional translators, taking into consideration their needs and preferences, whilst improving translators' MT literacy. This project is a starting point for facilitating communication ensuring that MT is developed and utilised effectively.

Central to our investigation is the role of vi-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹<https://inboxtranslation.com/resources/research/freelance-translator-survey-2023/> and <https://www.ciol.org.uk/ciol-insights-languages-professions-2>

sualisation systems in Explainable Artificial Intelligence (XAI) for Neural Machine Translation (NMT). We aim to assess the utility of various NMT visualisation tools for professional translators, examining how these tools contribute to their understanding of NMT models' decisions. Through a comprehensive review of existing explainability visualisation systems, we implement selected ones in the form of three demos available on HuggingFace Spaces, in order to determine their effectiveness in helping translators comprehend whether MT models produce accurate translations for appropriate reasons. By facilitating communication and understanding among key stakeholders, our project aims to promote the effective development and use of NMT systems.

The contributions of this paper are the following: (i) a revision and typology of state-of-the-art visualisation functionalities for the explainability for NMT; (ii) a set of ready-to-use explainability and visualisation tools available in the Hugging Face Spaces for the translator's use;² and, (iii) a translator-focused evaluation of explainability visualisations for NMT. The paper is structured as follows: Section 2 provides an overview of previous research on NMT explainability methods and existing visualisation systems. In Section 3, a typology of functionalities is examined. Section 4 details the methodology for assessing explainability functionalities for translators. Section 5 presents an analysis of the visualisation systems from the translator's perspective. Discussion and conclusions are presented in Sections 6 and 7.

2 Previous Research on Visualisation Tools for NMT

In this section, we present the state-of-the-art of visualisation methods and tools employed for Explainable NMT (XNMT). Visualisation is a key component of XNMT methods identified by Stahlberg (2020) in his survey paper. It is used in Model-intrinsic interpretability methods, post-hoc interpretability methods (interpreting predictions with input analysis), and the analysis of Confidence Estimation in Translation. In a more recent survey by Madsen et al. (2022), several post-hoc methods for NLP interpretability were reviewed, which has been further specialised by Leiter et al. (2023), who specifically reviewed methods for NMT metrics. In this work, we focus on the review

²<https://huggingface.co/gabrielanicole>

of existing systems that implement such methods, aiming to make XNMT accessible for translators.

2.1 XNMT Methods

We present Explainability Methods specifically for NMT implemented using a Transformer architecture (Vaswani et al., 2017). The translation process starts with an **input** sequence of words in the source language that undergoes different steps as defined in the **process**, and concludes with the generation of a **output** sequence of words in the target language (this process is detailed in Figure 1). We categorise XNMT methods into two types: *Inspection methods* and *Attribution methods*. Each type considers different aspects of the translation process to provide explanations to the final user.

2.1.1 Inspection methods

Inspection methods present a single-point decision made by the NMT system. The challenge lies in selecting valuable information directly from a model decision or parameter.

We inspect the NMT model in several parts of the process. On the input side, we consider the presentation of the tokenised input sequence that is being fed to the NMT system. On the output side, the presentation of the NMT probability of every generated token, and the visualisation of the decoding algorithm, such as the beam search sequence generation. These inspection methods are used as part of debugging techniques and provide transparency to the NLP pipeline (Alharbi et al., 2021). Inspection Methods can also be extended using manipulation procedures. In this case, the raw values of the NMT system are post-processed to be more easily interpretable. An example is the use of weights computed by the attention mechanisms to describe how the NMT system relates the source and output sentences (Wiegrefe and Pinter, 2019). The attention values have also been used to compute a Confidence estimation metric, as presented by Rikters and Fishel (2017).

2.1.2 Attribution methods

Attribution methods aim to elucidate the relationship between different parts of the translation process and the impact that one part has on another. These methods are often referred to as feature importance algorithms, as they model one part (e.g., the input tokens) as a set of features responsible for the generated output (Zhou et al., 2022).

There are different levels of attributions

(Kokhlikyan et al., 2020). *Primary* attribution focuses on the relationship between the input features and the corresponding generated outputs of the model (e.g. (Sundararajan et al., 2017), (Ding et al., 2019)). It uses the gradients (i.e., internal data) of the NMT with respect to the input, helping to visualise the impact of the input tokens on the output tokens. *Layer* Attribution variant extends attribution to all neurons in a hidden layer, and *Neuron* Attribution methods attribute specific internal, hidden neurons to the inputs or output of the model. For instance, Bau et al. (2018) presented a method to detect the neuron responsible for a particular linguistic property, and manipulating that neuron would alter the linguistic property in the output. Detecting the relationship between a specific part of the NMT model and a linguistic behaviour is also known as a *probing method*. For example, in Linguistic correlation Analysis (Dalvi et al., 2019a), a supervised method learns the most relevant neurons for an extrinsic task as Part-of-Speech classification. This approach helps uncover the linguistic properties encoded within the NMT model’s internal representations.

Inspection and Attribution methods can both be categorised as Model-Intrinsic or Post-Hoc methods, depending on whether they utilise the model’s internal data or only the inputs and outputs of the model. While Inspection Methods aim to explain a single decision made by the model, Attribution Methods are more complex as they analyse the interaction between different parts of the model that may not directly interact. Together, these methods provide both decision and model understanding of the NMT outputs.

2.2 XNMT Systems and Tools

While survey papers on explainability in AI encompass many systems from Computer Vision to Text Generation, existing reports on XAI (Phillips et al., 2021) or on Visual Analytics (e.g. (Cui, 2019)) do not focus on the task and processes of translation *per se*. Though acknowledging two main types of visualisation techniques for texts, Bodria et al.’s (2021) all-encompassing survey paper fails to capture all the investigation techniques based on visualisation that have been developed for NMT. Even if some NMT toolkits like THUMT (Tan et al., 2020) or JoeyNMT (Kreutzer et al., 2019a) propose cross-lingual attention as a standard functionality, visualisation is hardly exploited

to the best of its potential for XNMT. Instead, the focus is not exclusively on visualisations but rather on probing strategies (de Seyssel et al., 2022).

In the following, we review existing visualisation systems and subsequently propose a typology of implemented functionalities within these systems. This recap encompasses methods, toolboxes, or libraries used for visualising NMT.

2.2.1 Main Visualisation tools

Various tools are available that implement functionalities to explain the outputs and internals of NMT. These tools are available in the form of libraries and systems. Our analysis primarily focuses on visualisations designed for the Transformer architecture, but we also consider related tools that focus on sequence-to-sequence tasks (e.g. seq2seq-viz (Strobel et al., 2018)). The described tools offer a comprehensive overview of various functionalities that would enhance our understanding of translation as a task.

We review tools based on one of the following NMT toolkits: Fairseq (Ott et al., 2019), OpenNMT (Klein et al., 2017), JoeyNMT (Kreutzer et al., 2019b), and HuggingFace Transformers (HF-Transformers) (Wolf et al., 2020). The toolkit is the base of the XNMT tool, as the visualisation features are developed using the internals and outputs provided by the toolkit. We list XNMT method implementations detailed in Section 2.1. Table 1 summarizes eight analysed libraries and systems, detailing their creation year, NMT toolkit and if it is designed for the Transformer architecture (TR).

System	Year	NMT toolkit	TR
Seq2Seq-Vis	2018	OpenNMT	no
BertVis	2018	HF-Transformers	yes
Neurox	2019	HF-Transformers	yes
LIT	2020	HF-Transformers	yes
Captum	2020	HF-Transformers	yes
NMTViz	2021	py-torch	yes
Ecco	2021	HF-Transformers	yes
InSeq	2023	HF-Transformers	yes

Table 1: Libraries and systems overview. TR: Transformer.

2.2.2 Libraries

BertVis (Vig, 2019) is an inspection tool, which focuses on the NMT process by visualising the internals of the models, more specifically it provides detailed information of each multilayer and

multi-head attention of a neural model, supporting encoder-decoder architectures. It is specific for NLP models and works directly in Jupyter Notebook. *NeuroX* (Dalvi et al., 2019b) implements probing methods at a neuron and layer level. Additionally, it facilitates the manipulation of neuron values to explore architecture alternatives at tokenisation and neuronal levels. *NeuroX* also supports quality evaluation and analysis through Ablation studies, allowing users to analyse the impact of modifications on the generation of translated text, as demonstrated in previous research (Bau et al., 2018). *Captum* (Kokhlikyan et al., 2020) is a Python library designed for PyTorch models, offering access to model internals and computation of primary, neuron, and layer attribution methods. *Ecco* (Alammar, 2021) is an interactive inspection and attribution tool that operates within Jupyter Notebook. It supports a selection of models such as GPT-2, BERT, and RoBERTa. Similar to *Ecco*, *Inseq* (Sarti et al., 2023) is also based on *Captum* and provides comparable functionalities. However, *Inseq* extends its support to a wider range of models and is adaptable to various systems beyond Jupyter notebooks.

2.2.3 Systems

The following tools, presented in the form of systems, are standalone platforms tailored to facilitate explainability in NMT tasks.

Seq2Seq-Vis (Strobelt et al., 2018) is a system focused on aiding neural model developers in error detection through a set of inspection functionalities. It includes three main functionalities: (i) Inspection, presenting embedding space visualisation based on similar tokens from the training data for the encoder and decoder, attention visualisation between them, and probabilities for the generation of each token and the beam search; (ii) *What-if* Translations, allowing modification of selected tokens using the most probable ones, beam search, or attention between source and target tokens; and, (iii) Human error search for debugging, utilising the previous functionalities and relying on the NMT model’s understanding to identify bugs in the analysed architecture. It works on OpenNMT encoder-decoder architectures before the transformer era.

The *Language Interpretability Tool (LIT)* (Tenney et al., 2020) implements various tasks for explaining datasets, embeddings, and token representations. It utilises different primary attribution

methods for analysing model behaviour. Additionally, LIT visualises attention matrices, compares different data points and models, and provides performance evaluation. This versatile tool is compatible with various NMT toolkits, such as HF-transformer. *NMTVis* (Munz et al., 2021) Is the only tool that targets the translator user, offering functionality for exploring various translation alternatives using the generated target text and the beam search, along with the visualisation of attention between source and target sentences. The user can manually modify a translation, which updates the remainder of the target sentence, and navigate across different generation options to refine translations using the beam search.

While libraries are easier to incorporate into a new tool, systems are closed platforms that pose challenges when integrating with different models or third-party systems.

3 A Typology of Implemented Functionalities

In this section, we present a typology of implemented functionalities in state-of-the-art systems. To exemplify each functionality, we map them to the Transformer architecture (Figure 1, original figure taken from Vaswani et al. (2017)). Our survey of XNMT as a task and process follows an input, process, output analysis of the functionalities:

(i) Tokenisation (input/output) visualisation illustrates how input sequences are divided into tokens, which are then represented as embeddings in the encoder. Similarly, decoder output is presented in terms of tokens. Various XNMT tools display this functionality, like *BertViz* by showing attention links between tokens rather than words.

(ii) Embeddings (input/process) representation is depicted as a 2D or 3D projection through dimension reduction techniques, such as UMAP or t-SNE. As it is a space of points, multiple samples are used to relate different tokens. For instance, *LIT* illustrates the embedding space using several input sentences.

(iii) Attention weights (process) visualisation relates the input and output with its context at the encoder and decoder. It is represented as a bipartite graph (as in *BerViz*), or as a Heatmap Matrix (*InSeq*). In the encoder, self-attention relates the input sequence to itself. In the decoder, two attention types are used: self-attention relates the out-

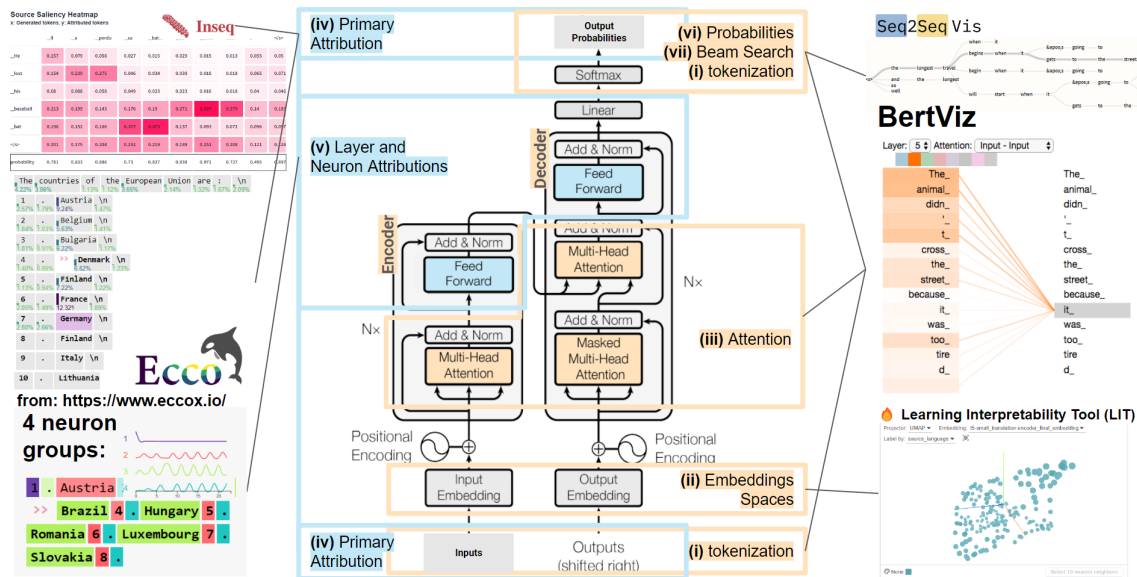


Figure 1: XNMT functionalities mapped in the transformer architecture and a corresponding example.

put to the generated tokens, and cross-attention relates the output tokens to the encoder output. This involves multiple layers and heads, each computing different attention weights, resulting in complex visualisations of multi-head and multi-layer weights for each attention type.³

(iv) Primary Attribution (Input/output) visualisation, which can be computed using different methods aims to illustrate the importance relationship between input and output tokens, attributing responsibility to specific input parts for generating each output part. While it reveals the relationship between input and output without explaining the process, it may utilise process information to compute a more accurate attribution metric. *Inseq* system presents this visualisation as a heatmap between inputs (rows) and outputs (columns), with darker colors indicating higher importance weight.

(v) Neuron and Layer Attribution (input/process/output) visualisation tries to pinpoint the responsibility of an output to a specific part of the architecture (a neuron or layer). For instance, *Ecco* illustrates this relationship by considering groups of neurons with a common linguistic task, and *NeuroX* use textual heatmaps to associate inputs with neuron values.

(vi) Probabilities (output) visualisation shows the prediction of the next token across the target

³Attention weights reveal internal model values, yet the link to model outputs is not clear (Jain and Wallace, 2019). We present this visualization to translators without assuming it offers explanations, enabling them to assess its utility.

dictionary. At each generation step, tokens more likely to appear next have higher probabilities. The visualisation displays top probable tokens, with darker colours indicating higher probabilities.

(vii) Decoding strategy (output) is the final step of the translation process. Here, the search for the sequence translation is exemplified, such as by visualising the beam search. This reveals each generated token step-by-step and how the optimal solution changes with respect to the search strategy and the beam size (i.e., the number of generated solutions). This visualisation is typically presented as a tree graph (e.g., *NMTVis*), offering the advantage of providing translation alternatives.

(viii) Training Data visualisation tries to present the datasets used to train the NMT systems. For example, *seq2seq-vis* uses them in the embedding representation to show similar input and output tokens, and *LIT* presents different data clusters and computes description metrics on them.

The described functionalities comprehensively map the Transformer architecture, although ongoing improvements are needed to develop better methods. In this work, we focus on evaluating them from the translator’s perspective.

4 Material and Methods

4.1 Data: Challenge sets

We adopted and, in part, adapted a challenge set (Isabelle et al., 2017), from which a selection of

example segments likely to be mistranslated by NMT were identified. Since the aforementioned challenge set has previously been used to examine translations into French, we decided to apply this to English-Spanish and English-Chinese tests. To assess the explainability of the visualisation tools, ten test segments were selected using a modified version of the challenge set created by Isabelle et al. (2017). Five segments come directly from Isabelle et al. (2017) challenge set and were selected for their varying morpho-syntactic properties, and five additional segments were created to complement this list, in Table 2 (details in appendix A).

(Isabelle et al., 2017)
1. The repeated calls from his mother [should] have alerted us.
2. The woman who [saw] a mouse in the corridor is charming.
3. I requested that families not [be] separated.
4. She was perfect tonight, [was she not]?
5. [Whom] is she going out [with] these days?
New test segments
1. The door [slammed shut].
2. He lost his [baseball bat].
3. The government’s new programme [was rolled out] last month.
4. [Berry] is a gifted student.
5. We will [leave no stone unturned] to hold [those responsible] to account

Table 2: Challenge sets

4.2 Models

We resorted to the Helsinki NLP opus models available on Hugging Face (Tiedemann and Thottingal, 2020). The three models (English-French, English-Spanish and English-Chinese) have an encoder-decoder Transformer architecture and use the Sentencepiece algorithm (Kudo and Richardson, 2018). They were chosen for pedagogical and interoperability purposes.

4.3 Visualisation

To present the functionalities to our translation experts, we developed an online web interface available on Hugging Face Spaces (details in Appendix B). Each functionality is built based on a specific state-of-the-art library, as follows:

Top-K and Beam Search Sequence: output

probabilities and decoding sequence generation based on *NMTVis*.

Attention: modified version of *BertViz* for the visualisation of attention weights.

Attribution: *Inseq* heatmaps of input X gradient method (Simonyan et al., 2013).

We explore how the explainability visualisations provide information about the challenge sets.

4.4 Explainability Evaluation

As a final global appraisal, we adapt the performance-explainability framework proposed by (Fauvel et al., 2020) to describe the translator analysis in specific for visualisation tools in XNMT. Following (Phillips et al., 2021), we include the evaluation of Meaningfulness, Accuracy, Knowledge Limits Explanations, as follows:

Meaningfulness: Is the explanation intelligible and understandable to the translator? Possible values: 1=no, 2=somewhat, 3=yes

Faithfulness: Can we trust the explanations? Possible values: the explanations are 1=incorrect, 2=imperfect, 3=perfect

Accuracy: Does the explanation accurately reflect the NMT processing? Possible values: 1=no, 2=somewhat, 3=yes

Knowledge limits: Does the explanation show the uncertainties of the NMT prediction? Possible values: 1=no, 2=somewhat, 3=yes

User: What is the target user category of the explanations? Possible values: 1=NMT expert, 2=translation expert, 3=broad audience

Usage: What is the intended use? Possible values: 1=debugging, 2=training, 3=professional use

Information: Which kind of information does the explanation provide? Possible values: 1=inspection, 2=inspection with post-processing, 3=attribution

We conducted a focus group with six translators working in English-Chinese (2), English-Spanish (1) and English-French (3). Each functionality is tested using the same ten source sentences by all users, and finally, the evaluation is the result of a group discussion with the support of NMT experts. We distinguish between translators and NMT experts because each group possesses a different set of skills and knowledge.

5 Results : XNMT evaluation

5.1 Top-K and Beam Search Sequence

Generally, the results produced by the Top-K and beam search sequence tool are both interesting and useful for viewing alternatives, particularly where synonymous words have been considered. For example, for English-Spanish, the example ‘I requested that families not be separated’ the final target translation uses the verb *solicitar* (literally, *to request*) however *pedir* (*to ask for/to request*) was also considered by the machine. Yet the information outlined within the Top-K feature demonstrates that the former received a higher probability, although the tokenisation of the word into three separate tokens *so-licit-é* (literally, *(I) requested*) makes it difficult to obtain any concrete statistical data confirming the probability of the word as a single lexical unit. Similarly for the English-French translation of this sentence, the Top-K visualisation shows a full list of synonymous translations for the English word *requested*, these being: *demande* (kept by the beam search algorithm), *exiger*, *prier*, *réclamer*, *vouloir*, *solliciter*.

In a similar light, when analysing the sentence ‘Berry is a gifted student’, which becomes *Berry es una estudiante talentoso* (literally, ‘Berry is a student talented’), and ‘Berry’ is tokenised as *Ber-ry*), it is difficult to assess the level of attention given to the more literal alternative such as *baya* (referring to the fruit as opposed to the name). In contrast, this type of tokenisation is useful for translating proper nouns, names in particular, into Chinese. Transliteration is the main way of addressing names between English and Chinese. The way in which ‘Berry’ is tokenised as *Ber-ry* shows how the model transliterates the name from a phonological perspective.

Moreover, within the same segment, the English-Chinese translation of the word *gifted* can be rendered in various ways, depending on the collocation embedded in the translation. The Top-K probable tokens can thus inform translators of the different possibilities available, therefore aiding translators to make more contextually-aware decisions. In addition, the complementary insights shown across the two visualisations are helpful for translators to navigate themselves among the different possible translations ranging from the level of tokens (i.e., within the Top-K) to the level of semantic trunks or even sentences (i.e., within the beam search sequence tree).

5.2 Attention

The visualisation tool of the multi-layer and multi-headed attention mechanisms can be instrumental in facilitating collaborations between computer scientists and translation practitioners in order to identify at what stage things go right or wrong during the processing stage, thus facilitating the potential to improve the overall performance of the MT model. However, its usefulness for translators is somewhat less optimistic. Within English-Spanish translation, the sentence ‘I requested that families not be separated’ yields interesting results whereby the subjunctive mood is correctly triggered due to a change of subject, yet the cross attention also demonstrates a high level of attention between the verb *se separaran* (literally, *they would be separated*) in the target translation and the subject of the verb in the English source text. Using the visualisation tool, it is possible to see that the particle *se* (a marker of the medial passive) places a greater amount of importance on the verb ‘separated’ – potentially suggesting the machine’s recognition of the English passive as a grammatical structure. Meanwhile, the verb *separaran* is tokenised as *separar-an*, with *-an* (the element indicating subject-verb conjugation) placing greater importance on the subject of the verb (families), which again may suggest the machine’s ability to recognise verb-subject agreement.

When analysing the English–French translation of the sentence ‘He lost his baseball bat.’, it is interesting to notice that the encoder’s self-attention shows a higher attention weight between the tokens *bat* and *baseball*. This might suggest that the presence of the latter word helps to disambiguate the polysemous word *bat*, and obtain the correct translation ‘Il a perdu sa batte de baseball.’

For English–Chinese translation, this tool is particularly helpful to identify where things start to go wrong, especially when analysing the ‘Cross Attention’. Using this tool, users can walk through the layers and locate the layer in which the information started to go wrong. For example, in the sentence ‘We will leave no stone unturned to hold those responsible to account’, the translation output is “我们将不遗余力地追究责任者的责任” (literally, we will spare no efforts to hold 责任者 *zerenzhe* [responsible person] accountable). Here, *zerenzhe* is not commonly used in this context; however, it is a literal translation for ‘those responsible’, as indicated within the different lay-

ers of the decoder.

Overall, the different language groups involved within this study consider the attention tool the most complex. In many cases, understanding layers as the different stages within the translation is fairly easy to grasp and thus deploy within teaching-based scenarios. For future translators, attention weights are a way to revisit an onomasiological/semasiological approach. Students could be asked to identify the most relevant links where constituents are properly delimited with the attention weights. Conversely, they could use their constituent detection competence to characterise the division of labour for the different layers. Nevertheless, an important limitation should be pointed out: attention-weight visualisation on long sentences is more difficult. From a translation perspective, it would be useful to gain insights into how the model processes long and complex sentences. And, provided with context-sensitive NMT models, it would be interesting to analyse greater-than-sentence-level textual features. This will help to assess the model's reliability on contextual analysis, for example, overall coherence of the translation, consistency of proper nouns and issues with co-referentiality. However, the current visualisation output of long sentences is difficult to interpret and thus the data becomes less meaningful.

5.3 Attribution

The attribution heatmaps have the potential to provide useful insights for the translator, particularly in the case of the source saliency heatmap, which makes it possible to see how words within the source text influence the final target translation. Similarly, the target saliency heatmap focuses on how the previously translated words influence the determination of the following words. Both tools can potentially allow translators to evaluate the efficiency of an MT model from the perspective of contextual cohesion, as well as the model's performance in producing natural collocations. However, the current version of the visualisation contains less focused information than a translator would need. A more interactive user interface might be helpful to enhance the usability of this tool. For example, when demonstrating the English-Spanish sentence 'The repeated calls from his mother should have alerted us', the source saliency heatmap yields no noteworthy results; in the English-Chinese direction, the heatmap shows

correct syntactic attention, but due to the fact that it failed to provide a semantically correct translation ('calls' mistranslated as 呼吁 *huyu* (appeal)), it can result in translators' confusion: is the visualisation trustworthy whereas the actual problem might be a lack of training data? However, the target saliency heatmap, shows an increased amount of saliency being given to the verb *deberían* (they should) to confirm its translation of *habernos alertado* (literally, having alerted us), which in Spanish is typically formed with the use of a modal verb such as *deber* (to have to) as we see here. There was one scenario in which the target saliency heatmap proved largely redundant when considering the sentence 'The woman who [saw] a mouse in the corridor is charming' as no colours appeared within the heatmap itself.

5.4 Global appraisal of Functionalities

Among the translators who tested the tools, three of them (one for each language pair) provided a fine-grained evaluation of every visualisation tool in terms of the criteria evoked in section 4.4. However, the following criteria were not rated by the translators but by the NMT experts: Faithfulness, Information, Accuracy and User. The reason for this is that in order to rate how accurately and faithfully a visualisation tool reflects the NMT processes, a detailed understanding of its inner workings is necessary, therefore this information was provided by the experts in the field. Similarly, for the Information, an in-depth understanding of the tool and data processing is needed. As for the User, we consider the definition of a "user by design", predefined by the creators of the tools. These criteria can be considered as being objective, while the remaining can be considered subjective. The latter were assessed by the translators. All the results can be found in figure 2. It presents separately the objective criteria (upper-left) and the subjective criteria, which are presented by language pair (and hence by evaluator).

We remark that the tools that were found most useful (for debugging, training and professional use) were the beam-search tree and the Top-K probabilities visualisation. These two also rank highest in meaningfulness, which indicates that this tool speaks to the translators the most. They also seem most capable of showing the limits in the model. This is probably due to their capacity to show alternative translations. Every

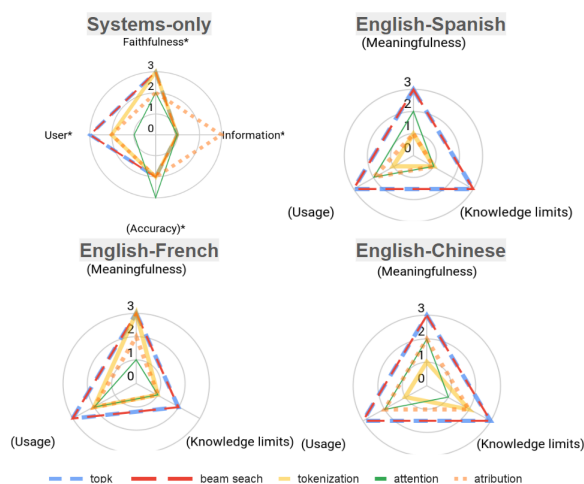


Figure 2: Global XNMT evaluation, upper-left graph shows the assessment done by NMT experts, followed by a graph per language pair (per evaluator) assessed by translators.

evaluator rated these two tools (Top-K and beam search) identically, which indicates a strong similarity between them. The attention visualisation tool proved to be only somewhat meaningful even though it is the most accurate of the visualisations. The evaluators agree that this visualisation doesn't permit them to detect the knowledge limits of the MT model, in fact this tool rated lowest in both Knowledge limits and Meaningfulness criteria. However, they do state that it can be useful not only for debugging purposes but as well during training of translation students. The visualisation of tokenisation is the tool on which the evaluators were less unanimous in terms of its meaningfulness. It was rated as not meaningful (1 out of 3) by two evaluators and as meaningful (3 out of 3) by one evaluator. This is probably caused by the fact that some evaluators found it misleading to see the probabilities of generating tokens rather than the the probability of the whole word. The Attribution method was rated mildly meaningful and not very capable of showing the model's limits. In combination with its low faithfulness and low accuracy due to the inner workings of this tool, it doesn't seem to be very useful for translators.

6 Discussion

6.1 XNMT for translation

From the translator perspective, one could hypothesise that there are parallels between traditional translation approaches (i.e., human translation) and the visualisation functionalities that we have tested. For example, Nida (1964) concept of

three-stage translation systems emphasises source text analysis, kernel extraction, the transfer, and the restructuring of meaning in the target language. The inspection methods of NMT (in particular, cross attention) resemble a deep understanding of the source text that is required to grasp its semantic and syntactic nuances; namely, the layers of attention forming part of the way in which NMT analyses and encodes the source sentence into representations capturing its meaning. Extracting the kernel, or comprehending the fundamental meaning of a sentence, could be linked to attribution methods (e.g., saliency), through which translators gain insights into the semantic and syntactic elements of a sentence that the model pays attention to. This process helps translators to understand how the model makes decisions and restructures the target translation to fit the norms of the target language. These links have the potential to help translators and trainees gain a superficial insight into the 'thinking process' of NMT models and expand their perceptions of the trustworthiness and reliability of such models. The inspection and attribution methods are similar to how a human translator might refine their understanding and approach to translation, which will lead to a continuous human-informed improvement cycle for NMT explainability.

6.2 Visualisation in CAT tools

We need to take into account the current computer-assisted translation (CAT) tools available to professional translators and discuss visualisation tools already at their disposal. Professional translators use Translation management systems (TMS), typically CAT tools, for their daily work. The main features of a TMS include project management, translation memory, terminology management, collaboration and review, reporting and analytics, automation and other systems integrations, etc. A TMS enables translators to leverage resources including translation memories, terminology databases and MT engines. This allows them to reuse previously translated segments and/or use raw MT output as a starting point for human translation, whilst maintaining consistency in terminology and phrasing. Since MT is usually an integrated feature of a TMS, translators either use MT suggestions as reference or directly post-edit the raw MT outputs. The working processes of a MT are not a primary concern for translators during

their regular workflows. Nevertheless, their ongoing automation anxieties need to be addressed and MT literacy is part of the overarching strategy.

It is important to take into account the levels of visualisation that translators and trainees can take in when promoting NMT explainability. The deep visualisations that we tested are very distant from their daily work. The visualisations in a TMS are functional-oriented which typically include the indicators of matches found in translation memories and terminology databases in the form of colour-coding, underlining, or other visual cues, the list of suggestions generated from concordance search, flagged potential quality issues, and progress bars, etc. These visualisations are set to present the complex information in a more intuitive and user-friendly manner, helping translators work more efficiently. In contrast, deep visualisation requires basic knowledge in NMT and clear guidelines to ensure correct interpretation. This is also the next step of our project: a workshop to disseminate the visualisation toolkit and to test the translators' and trainees' reception, and explore its wider usage.

6.3 Additional Functionalities

With the advent of Large Language Models (LLMs), it is tempting to use LLMs for Automatic Post-Editing of translations, a pipeline already implemented for Automatic Speech Recognition (ASR) systems such as WhisperingLlama (Radhakrishnan et al., 2023), which uses Llama (Touvron et al., 2023) to regularise and optimise Whisper's outputs (Radford et al., 2023) for ASR transcriptions.

We are already witnessing an integration of generative-AI with TMS; for example, the web-based system Wordscope. Along with the essential TMS features, the system integrates ChatGPT with ready-made prompts that allow translators to look up terms, search a topic or a concept, explore alternative expressions, back translate into the source language for quality check, proofread, post-editing, and more. A potential research question here is: can LLMs facilitate XNMT? Considering the increasing integration of generative-AI into the workflows of tech-savvy translators, is it possible to use LLMs to enhance the interpretation of NMT visualisations? For example, using generative-AI to help analyse the linguistic challenges that might be overlooked by human, and compare the results with the visualisations to fos-

ter more comprehensive evaluation.

In addition to the proposed research questions, the focus group highlighted several desires and requirements for XNMT to be fully deployed within the translator's workflow. One of the more divisive of which included potentially changing the presentation of subtokens within the final visualisation output (i.e., presenting the whole word as a single lexical item e.g., *berr-ry* as 'Berry'). Whilst it is generally understood that tokenisation forms an indispensable element of how the machine understands and process language (a feature enjoyed by the tech-savvy and developers of XNMT), the lack of a single overall probability for the entire lexical unit makes it challenging for translators to obtain meaningful statistical data that could be used to inform the translator's decision making processes.

7 Conclusion

In this paper, we have summarised the main visualisation tools adapted for XNMT, detailing the functionalities implemented in a prototype and discussing the potential benefits for translators. Our innovation lies in highlighting the translator's viewpoint and utilising XNMT to provide accountability for translators. This entails gaining a better grasp of the training data, monitoring the learning phase, or finding ways to understand the entire NMT process. As future work, we will continue exploring additional visualisation tools and evaluating their use, specifically focusing on one of the following translation moments: (i) initiating use of a new technology to understand NMT system workings during training, (ii) beginning a project to comparing, trusting, and selecting the best NMT, (iii) analyzing the translation process to identify reasons for poor output, and (iv) evaluating translation results e.g. to test alternatives.

Acknowledgement

This paper emanated from research supported by the MAKE-NMTVIZ project, funded under the 2022 Grenoble-Swansea Centre for AI Call for Proposals/ GoSCAI - Grenoble-Swansea Joint Centre in Human Centred AI and Data Systems (MIAI@Grenoble Alpes (ANR-19-P3IA-0003)). This work was also supported by the CREMA project (Coreference REsolution into MACHine translation) funded by the French National Research Agency (ANR), contract number ANR-21-CE23-0021-01.

References

- Alammar, J. 2021. Ecco: An open source library for the explainability of transformer language models. In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing: System demonstrations*, pages 249–257.
- Alharbi, Mohammad, Matthew Roach, Tom Cheesman, and Robert S Laramee. 2021. Vnlp: Visible natural language processing. *Information Visualization*, 20(4):245–262.
- Bau, Anthony, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James Glass. 2018. Identifying and controlling important neurons in neural machine translation. *arXiv preprint arXiv:1811.01157*.
- Bodria, Francesco, Fosca Giannotti, Riccardo Guidotti, Francesca Naretto, Dino Pedreschi, and Salvatore Rinzivillo. 2021. Benchmarking and survey of explanation methods for black box models. *ArXiv*, abs/2102.13076.
- Cui, Wenqiang. 2019. Visual analytics: A comprehensive overview. *IEEE access*, 7:81555–81573.
- Dalvi, Fahim, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James Glass. 2019a. What is one grain of sand in the desert? analyzing individual neurons in deep nlp models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6309–6317.
- Dalvi, Fahim, Avery Nortonsmith, Anthony Bau, Yonatan Belinkov, Hassan Sajjad, Nadir Durrani, and James Glass. 2019b. Neurox: A toolkit for analyzing individual neurons in neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 9851–9852.
- de Seyssel, Maureen, Marvin Lavechin, Yossi Adi, Emmanuel Dupoux, and Guillaume Wisniewski. 2022. Probing phoneme, language and speaker information in unsupervised speech representations. *arXiv preprint arXiv:2203.16193*.
- Ding, Shuoyang, Hainan Xu, and Philipp Koehn. 2019. Saliency-driven word alignment interpretation for neural machine translation. *arXiv preprint arXiv:1906.10282*.
- Fauvel, Kevin, Véronique Masson, and Elisa Fromont. 2020. A performance-explainability framework to benchmark machine learning methods: application to multivariate time series classifiers. *arXiv preprint arXiv:2005.14501*.
- Guerberof-Arenas, Ana and Joss Moorkens. 2023. Ethics and machine translation: The end user perspective. In *Towards Responsible Machine Translation: Ethical and Legal Considerations in Machine Translation*, pages 113–133. Springer.
- Isabelle, Pierre, Colin Cherry, and George Foster. 2017. A challenge set approach to evaluating machine translation. In Palmer, Martha, Rebecca Hwa, and Sebastian Riedel, editors, *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2486–2496, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Jain, Sarthak and Byron C Wallace. 2019. Attention is not explanation. *arXiv preprint arXiv:1902.10186*.
- Klein, Guillaume, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander Rush. 2017. OpenNMT: Open-source toolkit for neural machine translation. In Bansal, Mohit and Heng Ji, editors, *Proceedings of ACL 2017, System Demonstrations*, pages 67–72, Vancouver, Canada, July. Association for Computational Linguistics.
- Kokhlikyan, Narine, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqu Yan, et al. 2020. Captum: A unified and generic model interpretability library for pytorch. *arXiv preprint arXiv:2009.07896*.
- Kreutzer, Julia, Jasmijn Bastings, and Stefan Riezler. 2019a. Joey nmt: A minimalist nmt toolkit for novices. *arXiv preprint arXiv:1907.12484*.
- Kreutzer, Julia, Jasmijn Bastings, and Stefan Riezler. 2019b. Joey NMT: A minimalist NMT toolkit for novices. In Padó, Sebastian and Ruihong Huang, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations*, pages 109–114, Hong Kong, China, November. Association for Computational Linguistics.
- Kudo, Taku and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*.
- Leiter, Christoph, Piyawat Lertvittayakumjorn, Marina Fomicheva, Wei Zhao, Yang Gao, and Stefan Eger. 2023. Towards explainable evaluation metrics for machine translation. *arXiv preprint arXiv:2306.13041*.
- Madsen, Andreas, Siva Reddy, and Sarath Chandar. 2022. Post-hoc interpretability for neural nlp: A survey. *ACM Computing Surveys*, 55(8):1–42.
- Munz, Tanja, Dirk Vāth, Paul Kuznecov, Ngoc Thang Vu, and Daniel Weiskopf. 2021. Visualization-based improvement of neural machine translation. *Computers Graphics*.
- Nida, Eugene Albert. 1964. *Toward a science of translating: with special reference to principles and procedures involved in Bible translating*. Brill Archive.

- Ott, Myle, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*.
- Phillips, P Jonathon, P Jonathon Phillips, Carina A Hahn, Peter C Fontana, Amy N Yates, Kristen Greene, David A Broniatowski, and Mark A Przybocki. 2021. Four principles of explainable artificial intelligence.
- Radford, Alec, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR.
- Radhakrishnan, Srijith, Chao-Han Huck Yang, Sumeer Ahmad Khan, Rohit Kumar, Narsis A Kiani, David Gomez-Cabrero, and Jesper N Tegner. 2023. Whispering llama: A cross-modal generative error correction framework for speech recognition. *arXiv preprint arXiv:2310.06434*.
- Rikters, Matīss and Mark Fishel. 2017. Confidence through attention. In Kurohashi, Sadao and Pascale Fung, editors, *Proceedings of Machine Translation Summit XVI: Research Track*, pages 299–311, Nagoya Japan, September 18 – September 22.
- Sarti, Gabriele, Nils Feldhus, Ludwig Sickert, Oskar Van Der Wal, Malvina Nissim, and Arianna Bisazza. 2023. Inseq: An interpretability toolkit for sequence generation models. *arXiv preprint arXiv:2302.13942*.
- Simonyan, Karen, Andrea Vedaldi, and Andrew Zisserman. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*.
- Stahlberg, Felix. 2020. Neural machine translation: A review. *Journal of Artificial Intelligence Research*, 69:343–418.
- Strobel, Hendrik, Sebastian Gehrmann, Michael Behrisch, Adam Perer, Hanspeter Pfister, and Alexander M Rush. 2018. Seq2seq-vis: A visual debugging tool for sequence-to-sequence models. *IEEE transactions on visualization and computer graphics*, 25(1):353–363.
- Sundararajan, Mukund, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR.
- Talmy, Leonard. 1985. Lexicalization patterns: Semantic structure in lexical forms. *Language typology and syntactic description 3: Grammatical categories and the lexicon*, pages 57–149.
- Tan, Zhixing, Jiacheng Zhang, Xuancheng Huang, Gang Chen, Shuo Wang, Maosong Sun, Huanbo Luan, and Yang Liu. 2020. THUMT: An open-source toolkit for neural machine translation. In *Proceedings of the 14th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 116–122, Virtual, October. Association for Machine Translation in the Americas.
- Tenney, Ian, James Wexler, Jasmijn Bastings, Tolga Bolukbasi, Andy Coenen, Sebastian Gehrmann, Ellen Jiang, Mahima Pushkarna, Carey Radebaugh, Emily Reif, et al. 2020. The language interpretability tool: Extensible, interactive visualizations and analysis for nlp models. *arXiv preprint arXiv:2008.05122*.
- Tiedemann, Jörg and Santhosh Thottingal. 2020. Opus-mt–building open translation services for the world. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 479–480.
- Touvron, Hugo, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Vieira, Lucas Nunes, Carol O’Sullivan, Xiaochun Zhang, and Minako O’Hagan. 2023. Machine translation in society: insights from uk users. *Language Resources and Evaluation*, 57(2):893–914.
- Vig, Jesse. 2019. Bertviz: A tool for visualizing multi-head self-attention in the bert model. In *ICLR workshop: Debugging machine learning models*, volume 3.
- Wiegrefe, Sarah and Yuval Pinter. 2019. Attention is not not explanation. *arXiv preprint arXiv:1908.04626*.
- Wolf, Thomas, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Liu, Qun and David Schlangen, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October. Association for Computational Linguistics.
- Zhou, Yilun, Serena Booth, Marco Tulio Ribeiro, and Julie Shah. 2022. Do feature attribution methods correctly attribute features? In *Proceedings of*

A Appendix - Challenge Sets

Tables 3 and 4 present the full set of challenging translations used to evaluate each explainability visualisation tool. Following Isabelle’s 2017 procedure, we define each translation example with its related challenge.

These segments have been included within this study for their various grammatical and linguistic features. Firstly, (Talmy, 1985) distinguishes between two main language groups: those that favour conflation of path with motion (e.g., the Romance languages), and those that favour conflation of manner with motion (e.g., Germanic, and Slavic languages). Chinese appears to fall somewhere between the two, albeit with a slight preference towards the latter. Verbs of manner and path, which (Isabelle et al., 2017) call ‘crossing movement verbs’, present a difficulty due to the lexical and syntactic challenges arising when translating between languages that conflate information differently. The example sentence ‘The door [slammed shut]’ was used to examine how the model responds to such verbs. We also presented challenges involving prepositional verbs, which can result in inaccuracies concerning active vs passive voice, as well as syntax, or overly literal translations (i.e., within the sentence ‘The government’s new programme [was rolled out] last month’).

A second key difference between languages concerns the productivity of compounding as a means of word formation. Both English and Chinese are especially productive in this regard; however, this is not necessarily the case within languages such as French and Spanish. Compound nouns present difficulties by way of differences in phrasal word order (i.e., modifier + noun vs noun + modifier) in addition to potential issues with lexical ambiguity or polysemy. We used the sentence ‘He lost his [baseball bat]’ to test the model’s ability to identify polysemous words such as ‘bat’ (i.e., object vs animal). And finally, in addition to testing polysemy with compound nouns, we also tested polysemy within proper nouns or names, with a particular focus on examining the visualisation of data where names are likely to be transliterated. In this instance, the sentence ‘[Berry] is a gifted student’ was used.

Sentence	Challenge
The repeated calls from his mother [should] have alerted us.	Is subject-verb agreement correct? (Possible interference from distractors between the subject’s head and the verb).
The woman who [saw] a mouse in the corridor is charming.	Are the agreement marks of the flagged participles the correct ones? (Past participle placed after auxiliary AVOIR agrees with verb object iff object precedes auxiliary. Otherwise participle is in masculine singular form).
I requested that families not [be] separated.	Is the flagged verb in the correct mood? (Certain triggering verbs, adjectives or subordinate conjunctions, induce the subjunctive mood in the subordinate clause that they govern).
She was perfect tonight, [was she not]?	Is the English “tag question” element correctly rendered in the translation?
[Whom] is she going out [with] these days?	Is the dangling preposition of the English sentence correctly placed in the French translation?

Table 3: (Isabelle et al., 2017) challenges used to evaluate the explainability visualisation tools

B Appendix - Visualisation tools

In this section, we describe the implemented functionalities⁴.

General Interface We have created three demos available as spaces on the HuggingFace platform, all built using Gradio and Javascript. In Figure 3, we present the general interface, where translators can either choose a challenge or input a source text in English. The text is subsequently translated based on the selected model (en-zh for Chinese, en-es for Spanish, and en-fr for French).

Probabilities: Top-k Figure 4 shows the top-k most probable tokens to be generated, where in this case, k=10. The probability is represented on a scale of grey colours. At each generation step,

⁴Publicly available at anonymous

Translation

If challenge is selected from the challenge set list below

source text

target text

Challenge

category_minor

category_major

Challenge selection:

en-zh
 en-es
 en-fr
 en-sw

Figure 3: General Interface

Sentence	Challenge
The door [slammed shut].	Verb of manner and path - How has the manner and path been conflated, and does this follow the typical patterns of the target language?
He lost his [baseball bat].	Has baseball bat been translated as a compound noun, or two separate lexical items?
The government's new programme [was rolled out] last month.	Similar to verb of manner and path with added syntactical difficulties and passive vs active voice.
[Berry] is a gifted student.	Has Berry's name been translated literally? Transliterated?
We will [leave no stone unturned] to hold [those responsible] to account.	How has the idiomatic expression been translated? Has the syntax been adjusted accordingly?

Table 4: New challenge set used to evaluate the explainability visualisation tools

the top-k probable tokens are presented. According to the tokenisation used, one or several tokens could correspond to a single word. For instance, the word *alerter* was generated in two steps: first *alerte*, and then *r*.

Decoding Strategy: Beam Search Sequence Generation

The beam search visualisation is a simplified representation of the “beam search” decoding strategy, aiming to find the best “global” translation, i.e., the best sequence of translated tokens. Figure 5 displays the beam search decoding sequence generation using a beam size of 4. This visualisation presents the sequences (4) of output tokens in a tree structure, allowing users to notice the differences between alternatives. The top branch represents the sequence with the highest probability, while less likely sequences are displayed below.

Attention The attention visualisation shows the multi-layer and multi-head attention mechanism used in the transformer architecture. Each layer comprises several heads, each learning different weights between compared elements (tokens of the source or translated sentence). In the visualisation, each head is represented by a colour, with darker colours indicating higher attention weights. This information is represented through connection lines and coloured boxes. Three attention options are presented in the visualisation: (i) *encoder self-attention*, which relates the tokens of the source text to each other; (ii) *decoder self-attention*, which relates the translated tokens to the previously generated tokens; and (iii) *cross-attention*, which relates the translated tokens and

Exploring top-k probable tokens

Les	appels	répété	s	de	sa	mère	auraient	dû	nous	alerte	r	:	</s>	</s>
Ses	rappel	réitéré	de	lancés	la	maman	[auraient:0.45271835]	avertir	nt	!	-	</s>
L	cris	de	es	venant	Sa	Mère	devraient]]	prévenir	z	</s>	-	-
La	coups	répét	des	que	son	mé	devraient	été]	a	s	!	C	!
Ces	nombreux	répétées	S	qu	cette	femme	ont	eu	vous	appeler	ment	.	"	!
Il	multiples	récurrent	,	émanant	ses	mer	aurait	aurait	les	alarme	ra	...	Je	,
"	demandes	à	d	[leur	mères	doivent	[avoir	averti	.	"	"	...
Le	conversations	que	et	par	ma	m	(fallu	m	sensibiliser	rs	[[."
C	appel	multiples	.	,	ta	père	,	fait	s	attirer	R	:	([
Des	répétition	successifs	par	provenant	notre	famille	nous	auraient	le	alerte	ner	de	Il	:

Figure 4: TopK probable tokens

Exploring the Beam Search sequence generation

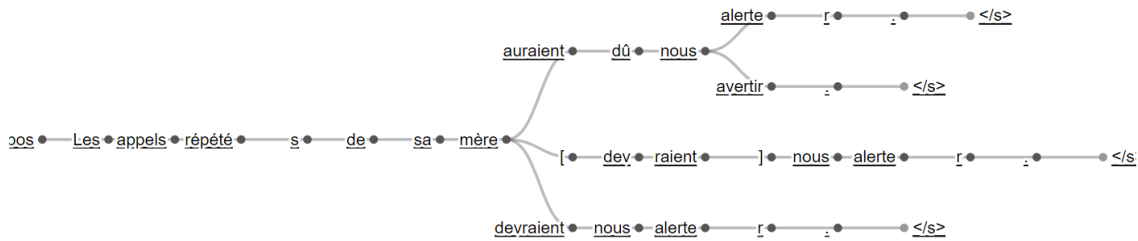


Figure 5: Beam Search Sequence Generation

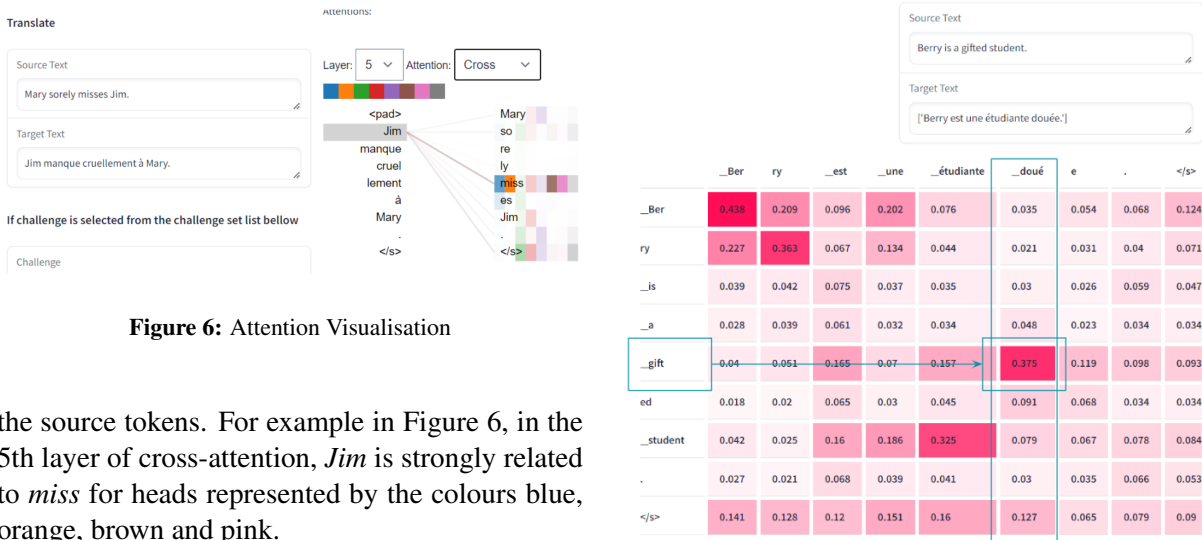


Figure 7: Primary Attribution

the source tokens. For example in Figure 6, in the 5th layer of cross-attention, *Jim* is strongly related to *miss* for heads represented by the colours blue, orange, brown and pink.

Attribution The attribution visualisation presents the importance of each token of the source text (rows) in generating the tokens of the translated text (columns). This attribution is computed using the input X gradient method (Simonyan et al., 2013). In the heatmap, the importance of compared tokens is indicated by the darkness of the colour. For example, in Figure 7, the most important token for generating *doué* is *gift*.

TopK and Beam Search Sequence Generation functionalities are based on state-of-the-art tools. However, they are implemented by us. For attention visualisation, we adapted the *BertViz* library to make it compatible with Gradio, while the *Inseq*

library made possible the attribution visualisation.

Mitigating Translationese with GPT-4: Strategies and Performance

Maria Kunilovskaya¹, Koel Dutta Chowdhury¹, Heike Przybyl¹,
Cristina España-Bonet², and Josef van Genabith^{1,2}

¹Saarland University, Saarland Informatics Campus, Germany

²German Research Center for Artificial Intelligence (DFKI)

maria.kunilovskaya@uni-saarland.de

Abstract

Translations differ in systematic ways from texts originally authored in the same language. These differences, collectively known as translationese, can pose challenges in cross-lingual natural language processing: models trained or tested on translated input might struggle when presented with non-translated language. Translationese mitigation can alleviate this problem. This study investigates the generative capacities of GPT-4 to reduce translationese in human-translated texts. The task is framed as a rewriting process aimed at modified translations indistinguishable from the original text in the target language. Our focus is on prompt engineering that tests the utility of linguistic knowledge as part of the instruction for GPT-4. Through a series of prompt design experiments, we show that GPT-4-generated revisions are more similar to originals in the target language when the prompts incorporate specific linguistic instructions instead of relying solely on the model’s internal knowledge. Furthermore, we release the segment-aligned bidirectional German–English data built from the Europarl corpus that underpins this study.

1 Introduction

There has been a surge of interest in the impact of translationese on the performance of natural language processing (NLP) applications. Translationese has been shown to have tangible effects on

the outcomes of various cross-lingual tasks, potentially leading to biased results and decreased or artificially inflated performance, especially in evaluating machine translation (MT) models (Zhang and Toral, 2019; Graham et al., 2020), but also in the natural language inference tasks when using translated datasets and cross-lingual transfer scenarios (Artetxe et al., 2020). While translationese is viewed as an inalienable property of translated language, preferences may lean toward translation variants that are closer to target language patterns provided that the meaning and usefulness of the message in the source language (SL) are retained. The task of reducing translationese by making translations less deviant from the originally authored text in the target language (TL) is a newly recognised and relevant NLP problem. At the same time, only a few studies actively address it, including Dutta Chowdhury et al. (2022) who remove translation bias in latent representation space, as well as Jalota et al. (2023) and Wein and Schneider (2024), debiasing translations at the surface text level.

Our work is the first to explore the utility of linguistically informed prompts to harness the generative capabilities of large language models (LLMs) in the task of translationese mitigation. This approach is inspired by the successful application of LLMs to a range of text adaptation tasks including simplification (Feng et al., 2023), style transfer (Suzgun et al., 2022; Reif et al., 2022), and translation (post-)editing (Chen et al., 2023; Rana et al., 2023). To the best of our knowledge, only Chen et al. (2023) uses LLMs to address translationese reduction. We extend this line of research.

Specifically, we focus on exploring the impact of linguistic knowledge, made available to

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

the LLM via prompts, on the outcomes of translationese reduction. The key research question is **what type of information is required in the prompts to effectively guide the model through the rewriting process**. We propose two approaches: (i) a self-guided approach, which probes the ability of the LLM to solve the task independently using its internal knowledge versus (ii) a feature-guided approach, which relies on detailed linguistically-informed instructions to edit the input. The instructions are based on the observed deviations of each individual segment from the expected TL norm. We define the expected TL norm as the type of language that can be expected in the target culture in a comparable communicative situation. It is represented by the average feature values from the register-comparable corpus of TL documents produced by native speakers of the TL (hereinafter referred to as originals).

The contributions of this work are as follows:

- We formulate the translationese mitigation task in an LLM-prompting setup, where an LLM is expected to remove the translation process artefacts and generate a ‘less translated’ version for an existing human translation (HT).
- We demonstrate the importance of detailed linguistically-informed instructions in formulating prompts, individually tailored for each segment.
- We release the document- and segment-level aligned corpus created from Europarl for this study and the multiparallel datasets for English–German and German–English contrastive samples including LLM generated versions aligned with the inputs¹.

These contributions collectively address our research question and advance our understanding of the impact of linguistic knowledge available to the LLM via prompts on the outcomes of translationese reduction. The remainder of this paper is organised as follows: Section 2 discusses related work. In Section 3, we introduce our prompt generation approaches. Section 4 details our experimental settings, including the rationale behind our linguistic feature design, feature extraction and selection methods, data description and our evaluation strategy. Section 5 presents and discusses the results. We conclude with a summary in Section 6.

¹<https://github.com/SFB1102/b7-b6-prompting-eamt2024>

2 Related Work

Translationese artefacts exert a substantial influence on diverse downstream tasks. In MT, Toral et al. (2018) and Edunov et al. (2020) found that source sentences that were already the result of a translation were easier to translate than non-translated sources returning higher BLEU scores. Graham et al. (2020) and Zhang and Toral (2019) also showed that translationese in test sets could lead to inflated and inaccurate evaluation scores and recommended non-translated sources in MT evaluation to avoid these biases. The influence of translationese on MT goes beyond evaluation. For example, Riley et al. (2020) trained the translationese classifier to tag the sentences in training data to control the output domain: translationese (“Tr”) or original/natural text (“Nt.”). In other cross-lingual applications, Singh et al. (2019) showed that substituting original training samples with their translations from another language improves performance on natural language inference tasks. Clark et al. (2020) introduced a translation-free question-answering dataset to avoid having inflated gains from translation artefacts in transfer-learning tasks. Artetxe et al. (2019) found that cross-lingual models suffered from induced translation artefacts when evaluated on translated test sets.

Active attempts to level out translationese bias include a method that can be applied in the *translate-train*² cross-lingual setup (Yu et al., 2022). They created a mapping from the original to the translated language, projecting original and translated text into a shared multilingual embedding space and minimising the distance between the mapped representations of the originals and translations. To mitigate translationese effects in translated data, Dutta Chowdhury et al. (2022) extended the Iterative Null Space Projection algorithm (Ravfogel et al., 2020) originally designed to mitigate gender attributes, to *debias* translationese artefacts, and not directly on the text itself, which makes them less interpretable. Wein and Schneider (2024) reduced translationese deviations at the surface level of text using Abstract Meaning Representation (AMR) proposed by Banarescu et al. (2013) as an intermediate form to abstract away from translationese artefacts. In another line of research, Jalota et al. (2023) reframed the

²In this setting, the training is based on translated data instead of originally authored data.

task as a self-supervised monolingual translation-based style transfer task, aiming to make human-translated text closely resemble original texts in the TL. However, whether current out-of-the-box LLMs are able to mitigate translationese from text without removing traces of other variables remains unexplored. Apart from the key related works in translationese mitigation, we elaborate on other contemporary studies that have used LLMs for manipulating text, sometimes with goals related to refining translations or removing undesired information from text representations. Vilar et al. (2023) benchmarked the capabilities of LLMs to translate, and Kocmi and Federmann (2023) and Lu et al. (2023) to evaluate translations. Along the same line, Hendy et al. (2023) extensively analyses the translation output of LLMs to demonstrate that GPT-enabled translation achieves high quality when utilised for the translation of high-resource languages. However, it still falls short in terms of translation quality for underrepresented languages. Likewise, Raunak et al. (2023) investigated these differences in terms of the literalness of translations produced by standard NMT and ChatGPT-3.

Contemporaneously to the present work, Chen et al. (2023) propose a simple way to refine translations iteratively with LLMs based on automatic post-editing that imitates human corrections.

3 Prompt Generation

Our experiments are designed to explore the effectiveness of including various types of information in prompts that influence the generative behaviour of an LLM in the task of translationese mitigation. The study is based on a bidirectional German-English subset of Europarl data. Each translation direction is aligned at the segment level, meaning that depending on the syntactic arrangement of the same content the source or the target side of the parallel data can have more than one sentence. We experimented with two prompting approaches: self-guided and feature-guided, each with two modes (min and detailed). The full prompt examples for each of these four prompting setups appear in Appendix C. The four setups vary in the degree of independence in decision-making given to the model and in the level of linguistic instruction. Below we provide a description for each setup.

1. **Self-guided modes:** These modes rely on the model’s discretion in solving the task.

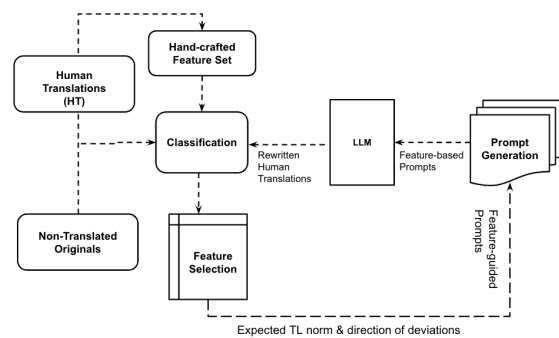


Figure 1: An overview of our pipeline based on the feature-guided approach.

min: In this mode, the prompt formulates the translationese reduction task without any reference to the concept of translationese or any other linguistic knowledge, in layman’s language: *Your task is to re-write a human translation in a more natural way if necessary.* Importantly, the model is given the option to return the input if it does not detect any traces of translationese, i.e. if the translation already sounds *like a text originally produced in the target language.*

detailed: Unlike the previous setup, the prompt contains a concise paragraph (186 words) explaining the concept of translationese as discussed in translation studies (Volansky et al., 2015; Hu and Kübler, 2021). It describes the known trends in translator behaviour and typical translationese indicators established in the literature. The option to return the input translation in case the model could not detect translationese deviation is kept. Figure 1 shows an overview of our pipeline for the feature-guided approach.

2. **Feature-guided modes:** The prompts include specific linguistic instructions that limit the model to a set of required transformations for each input translation. The list of instructions is tailored for each segment and addresses the most prominent deviations of this segment from the expected TL norm based on a number of linguistically motivated hand-crafted features (Section 4.1). The TL norm for each feature is calculated as the average across all segments in the original text category. The instructions for a particular feature are included in the prompt if the feature met the following criteria: (i) deviated more than

2.5 times from the TL norm in the direction observed in HT (e.g. in German translations, the frequency of additive connectives was lower than in non-translations, while translations into English had significantly more additive connectives than comparable non-translations in English), (ii) it was among the top 15 translationese indicators as flagged by SVM feature weights for each translation direction. If none of the features exceeded the 2.5-ratio threshold, the segment was not sent to the model and remained unchanged. The instructions for all segments were pre-compiled based on the threshold calculations and formatted as a newline-separated list appended to the task statement, source segment, and target segment (i.e. HT). The two variations of this setup were only different in how detailed the description of each instruction was.

min: The model was given a task *to re-write a human translation in a more natural way* by following the pre-compiled instructions. The instructions were formulated in a very concise manner. For example, *Make causative-consecutive relations between parts of the sentence more explicit.*

detailed: The task and the instructions were explained in more detail, offering descriptions of the linguistic concepts. Where possible, we provided lists of TL-specific examples for linguistic categories. Those prompts started with a brief definition of translationese followed by instructions like *Make causative relations between parts of the sentence more explicit. This can be done by using connectives like: because, therefore, so that, for this reason, as a result, after all, for that reason, hence, consequently, to this end.* In formulating the descriptions we relied on the definitions from the UD framework.³

In summary, in the two self-guided modes, the LLM’s behaviour is not constrained by specific rewriting instructions. The model had to make self-guided decisions not only on how to rewrite a segment but also on whether any transformation was necessary at all. In contrast, the two feature-guided modes closely supervised the model by

³<https://universaldependencies.org/u/dep/index.html>

specifying linguistic properties to be transformed in the rewriting process. All prompt types contained the source segment and its human translation. Preliminary experiments indicated that when the model was not constrained by the source segment, the re-writing process was highly volatile. Throughout this study, we only considered segments longer than eight words.

LLM Specifications. For our experiments, we use the GPT-4 model through the OpenAI API.⁴ This model returned more consistent results than GPT-3.5-turbo in a preliminary study. Our best results are obtained with GPT-4 and the default temperature (0.7). Although we attempted to suppress noise⁵ in the GPT-4’s output by appending formatting instructions to each prompt (e.g. *Do not add any meta-phrases or quotation marks*), the rewritten versions required extensive cleaning. The model’s comments were varied and the output had to be manually curated. Interestingly, even though the instructions were provided in English, the model added meta-comments either in German or in English when working on re-writing translations into German.

4 Experimental Settings

4.1 Linguistic Features

We propose to capture translationese with a set of morpho-syntactic features and text measures extracted from the Universal Dependencies (UD) annotation of the data. Unlike surface features like ngrams and neural network-based feature-learning approaches to translation detection, explicit discrete structural features have a lower risk of capturing irrelevant topical differences between the categories (Volansky et al., 2015; Borah et al., 2023). They are more interpretable and can be incorporated into human-readable rewriting instructions for an LLM. The initial feature set included 58 features and was motivated by previous research in language-pair-specific translationese (Evert and Neumann, 2017; Kuniilovskaya and Lapshinova-Koltunski, 2020) and contrastive studies (Konig and Gast, 2007), as well as multilingual analysis (Hu and Kübler, 2021). In Appendix A, the fea-

⁴<https://platform.openai.com/docs/guides/gpt>. The final version of the re-written translations analysed here was obtained between 08 and 10 March 2024.

⁵refers to undesirable outputs in model-generated text, including unwanted copies of the input, additional quotes and meta-comments from the model like: ‘Here is the revised translation.’

tures are categorised according to the type of linguistic units they capture. Our feature set contains grammatical forms, morphological word classes, clause types, syntactic dependencies, word order patterns, discourse elements, and textual measures. Generally, we gave preference to the features that:

- captured relatively frequent linguistic items to minimise sparsity as much as possible, especially at the segment level,
- were suggested as contrastive for the given language pair and/or were expected (or known) to generate translationese deviations from the TL norm.

Feature Extraction. For most features (37 out of 58), the extraction was straightforward and directly dependent on the accuracy of automatic annotation. The annotation quality is comparable across our languages, according to the official report for the models⁶ used. Six features of the remaining 21 features (various discourse marker types and adverbial quantifiers) relied on external pre-defined lists which were compiled using previous research in language variation for each language (Biber, 1988; Nini, 2015; Evert and Neumann, 2017), while the other 15 features included (i) straightforward metrics such as sentence length in tokens, word length, number of simple sentences, number of clauses per sentence, the ratio of core verbal arguments expressed by nouns, (ii) mean hierarchical distance and mean dependency distance (Jing and Liu, 2015), (iii) type-to-token ratio calculated as the ratio of part-of-speech-disambiguated content word types to their tokens, lexical density calculated as the ratio of disambiguated content word types to all tokens, (iv) and six word-order patterns that were discussed as English-German contrasts (Konig and Gast, 2007). All features were estimated and normalised at the sentence level and mean-aggregated for segments or documents. The highly correlated features were excluded (cutoff=0.65 for both languages).

Feature Evaluation and Importance. Table 2 shows that the proposed feature set demonstrated relatively high classification results at the document level. The feature selection did not yield considerable gains in performance: the improvements on the optimal 29 and 45 features (reported in Ta-

⁶<https://stanfordnlp.github.io/stanza/performance.html>

ble 2) were in the fractional part of the scores. This suggests that the proposed feature set does not include irrelevant features and is effective in capturing translationese. None of the features could reliably distinguish the categories on its own, demonstrating that translationese is a subtle phenomenon, which is better captured through feature patterns, in a multi-variate setup.

4.2 Data

We use the Europarl-UdS preprocessing pipeline⁷ to extract parliamentary speeches⁸ delivered in German and English by native speakers and their translations into English and German respectively. Our rewriting approach required parallel data, therefore, we report the details on sentence alignment quality. The documents were automatically aligned with LF Aligner⁹, a wrapper over the *hunalign* library (Varga et al., 2005), using domain-specific bilingual glossaries built from IATE dictionaries.¹⁰ The resulting parallel corpus was limited to the documents with an average document-level similarity score returned by the alignment tool over 0.3 and 0.5 for German-to-English and English-to-German directions, respectively. The manual evaluation of the automatic alignment, performed by a compensated research assistant on 80 document pairs (750 sentence pairs) randomly extracted for each direction, revealed that the resulting parallel corpus contained at most 4.5% (German-to-English) and 1.8% (English-to-German) of misaligned segments.

For this study, the corpus was balanced across translation directions by taking 1500 random document pairs that contained at least 450 tokens in the source language. The document length filter excluded short documents containing formulaic exchanges between the Chair and the participants of the debates in the European Parliament. All textual data were automatically parsed with the default Stanza packages for German and English (Qi et al., 2020). The quantitative parameters of the research data are given in Table 1.¹¹

⁷<https://github.com/chozelinek/europarl>

⁸It is well known that translation direction and register are the two major factors that influence the properties of translations (Redelinghuys, 2016; Evert and Neumann, 2017; Kunilovskaya and Lapshinova-Koltunski, 2020; Kunilovskaya and Pastor, 2021). Europarl data is convenient because it helps control for these factors.

⁹<https://sourceforge.net/projects/aligner/>

¹⁰<https://iate.europa.eu/search/standard>

¹¹The datasets are available as an indexed long table here: <https://zenodo.org/records/11127626>

		docs	segs	tokens
DE	original	1500	38,305	967,385
	translated	1500	36,078	924,919
EN	original	1500	36,078	927,045
	translated	1500	38,305	1,060,295

Table 1: Parameters of the entire research corpus (after filtering and annotation). EN (English) and DE (German) stand for the language of the comparable samples of originally authored and translated text. All translations are from the other language in the language pair. For example, DE translated are translations into German from English. DE original are texts in German by German native speakers.

The corpus in Table 1 was further distilled to obtain a contrastive sample of 200 documents in each TL that concentrated the translationese-related phenomena. To this end, we ran a 10-fold binary document-level translationese classifier using the features described in Section 4.1 and classification setup from Section 4.3. The results of this classification can be found in Table 2. For comparison, we report results for the full feature set and the optimal set of features (see details on feature selection in Section 4.3).

		feats	docs	F1
DE	29	3000	88.83 \pm 1.99	
	58		88.39 \pm 2.54	
EN	45	3000	80.05 \pm 1.68	
	58		79.66 \pm 2.05	

Table 2: The quality of the document-level translationese classifications across the two languages in the 10-fold cross-validation setup. The average document length in the translated text categories is around 700 tokens, 25.5 segments.

The contrastive subset was defined as 100 ‘most translated’ and 100 ‘most original’ documents based on the probability over 0.99 of belonging to their true class returned by the classifier on the best-performing 29 and 45 features for German and English, respectively. This data filtering step was required to meaningfully downsize the data to a subset manageable in the prompting experiments. Given the relatively high quality of the translationese classification (F1 score of 88% for German and 80% for English in Table 2), we have good reasons to believe that the selected documents bring into focus the contrasts between translations and non-translations while being naturally-occurring texts containing cohesive sequences of sentences. The parameters of this experimental subset appear in Table 3.

		segs	tokens	seg_len \pm std
DE	original	1908	59,942	31.4 \pm 17.6
	translated	1934	57,492	29.7 \pm 14.1
EN	original	1987	55,128	27.7 \pm 13.0
	translated	1919	65,065	33.9 \pm 19.6

Table 3: Parameters of the contrastive subset for rewriting experiments. Note that the originals here are not the sources for the translations in the other language. Instead, they are the top documents predicted as originals by the classifier (Table 2).

4.3 Evaluation

Translationese Classification. Our main translationese mitigation evaluation method is segment-level¹² text classification. If a rewriting strategy is effective, the accuracy scores for classifying translationese on the rewritten output should be lower compared to classification on HT (human-translated) text. In other words, there should be a negative difference in accuracy scores between the rewritten output and the initial HT, indicating that the rewritten versions blended better with the TL norm than the existing HT. For all experiments, we used a simple Support Vector Machine (SVM) with a linear kernel (C=1) in a 10-fold cross-validation setup. Linear SVM was preferred because it allows access to feature weights. The feature weights were used to identify a set of 15 most informative features. These features were used in prompt engineering and for evaluation purposes. The feature selection was performed using Recursive Feature Elimination technique with a linear SVM as implemented in the scikit-learn library.¹³ All classification results are reported for the top 15 features and for the full feature set. Although the number of instances per category was almost the same, we report a macro F1 score throughout to avoid any impact of the data imbalance on the results.

Re-translation (RT). As a sanity check for the rewriting approaches outlined in Section 3, we ran a re-translation mode (referred to as RT) to ensure that in the rewriting setups, the model follows our instructions and edits the existing translation, rather than returning a new translation. Here, we prompt the model to re-translate an existing HT if it detects any translationese deviations.

¹²Rewriting experiments on documents resulted in cropped GPT-4 output and therefore segment level was preferred.

¹³https://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.RFE.html

Statistical Analysis. The analysis of the classifiers’ performance was supported by tracking the shifts in the feature values observed in the generated text against original texts and HTs. This helped us understand whether the model managed to level out the existing translationese deviations and whether it introduced new tendencies. The significance of differences between originals in the TL and rewritings was estimated using the two-tailed Mann-Whitney U Test for independent samples. The results are considered significant at the confidence level of 5%.

Content Preservation. We evaluate the quality of the GPT-4 outputs in preserving the meaning of the input translations using COMET (Rei et al., 2022). We use two variants of COMET for this purpose: (a) R, reference-based (wmt22-comet-da) and (b) QE, the reference-free COMETQE (wmt20-comet-qe-da).

Manual Analysis. The automatically edited translations and re-translations were evaluated by one of the authors of this paper, a German-native professional translator with English and German as their working languages. The evaluator reviewed a random sample of 25 generated rewritten segments for each mode and translation direction. These segments were presented in the context of the source segment and the professional HT. Their task was to assess translation faithfulness to the source (accuracy) and lexicogrammatical acceptability (fluency) using a 1-6 scale (higher is better) for each output mode. Additionally, they checked whether the rewritten translations were compliant with the provided instructions (feature-guided modes only) to see whether the model followed the instructions. The expert was not asked to pass judgments about the translationese properties of the items in their sample. We maintain that translationese is a property of language that is visible to a machine rather than a human.

5 Results and Discussion

Translationese Classification. The results of our baseline SVM segment-level classification between originals and HTs from the contrastive sample (see Section 4.2) in each TL are reported in Table 4. We report F1 scores on the top 15 features and on the full feature set to throw the performance on the top 15 features into perspective.

The main observations from Table 4 are:

	feats	segs	F1
DE	15	3842	81.06±0.76
	58		81.51±1.79
EN	15	3906	75.60±1.87
	58		78.30±1.42

Table 4: Segment-level classification results on human translations from the contrastive 200-document sample using linear SVM. EN and DE stand for the target language.

(i) HTs into German contain more machine-detectable deviations from non-translations than translations into English, (ii) the reduced 15-feature set returns results comparable to the full 58-feature set, especially in German. We address these strong translationese predictors in the GPT4-based rewriting pipeline.

To assess the impact of rewriting on translated segments from the contrastive sample, we conduct another set of translationese classifications using the same original texts and their GPT4-rewritten versions on the top-15 subsets of translationese indicators addressed in the rewriting process and on the full-58 feature set. Table 5 shows the differences in F1 scores. Below we show some ob-

		Rewriting Setups				
		RT	Self-guided		Feature-guided	
		-	Min	Detail	Min	Detail
DE	15	-0.28	-0.27	-1.01	-2.39	-2.21
	58	0.10	0.53	-0.56	0.06	-0.28
EN	15	-3.32	-2.70	-4.10	-3.18	-7.63
	58	-0.58	-1.40	-1.61	-1.61	-4.07

Table 5: Differences in F1 scores between the segment-level results on the rewritings and on human translations from the contrastive sample (Table 4). The best results for each feature set are shown in bold.

servations from these results. Recall that lower translationese classification accuracy would suggest that rewritten segments became less distinguishable from originals after editing. The negative differences in Table 5 confirm that GPT-4 can be conditioned through prompting to address the task, even if the overall gains are small on the segments from the contrastive 200-documents sample. The rewriting task is more successful in English than in German. All attempted approaches decrease the prominence of translationese in the English translations by at least 0.58 points. In particular, when given detailed instructions based on the linguistic features (*Feature-guided Detail* mode), we observe a substantial 7.63 and 4.07 percentage

point decrease in classification results for the top-15 subset and for the full-58 feature set, respectively, working with the segments from the contrastive sample of 100 originals and 100 translations.

For German, the best-performing modes are the *Feature-guided Min* for 15 features and the *Self-guided Detail* setup for 58 features. Table 5 shows that the results are better for the 15 strong translationese predictors, specific for each language, even for GPT4 rewriting modes that did not rely on features. It means that the model effectively picked and reduced the most prominent translationese deviations even when it was not prompted to do so. The modes with the linguistic explanation (Detail) seem to be better than *Min* mode regardless of whether the model was presented with a list of specific rewriting instructions or was left to decide how to tackle this text adaptation task (except the *Feature-guided* approach for German on 15 features). Feature-guided modes were on average more successful than self-guided modes, especially for English. The performance on the features that were addressed in the instructions shows that the instructions were carried out in the rewriting.

Finally, the comparison with the SVM classification outcomes for the re-translation task indicate that the model did not simply return a new translation of the source. Although the model reduced translationese in the re-translation task, the explicit editing tasks performed better (cf. RT column to Detail columns in Table 5). Overall, the properties of rewritten documents are shifted towards being more similar to original texts. This effect is visible in Figure 2 which displays Kernel Density Estimation (KDE) plots for the values on the ‘translationese’ component obtained through Principal Component Analysis (PCA) of the 58-dimensional feature space. These plots capture the distribution of the values on this PCA component for original, human-translated and LLM-rewritten segments. Figure 2 shows that the rewritten documents (in red) are shifted from the area taken by the translated texts on the right (green line) to the non-translations’ left side of the graph.

Statistical Analyses. First, we find an imbalance in the segments affected by the diverse rewriting approaches across TLs. Recall that in the self-guided and re-translation modes the model was given the option to return the input unmodified while in the feature-guided modes, the seg-

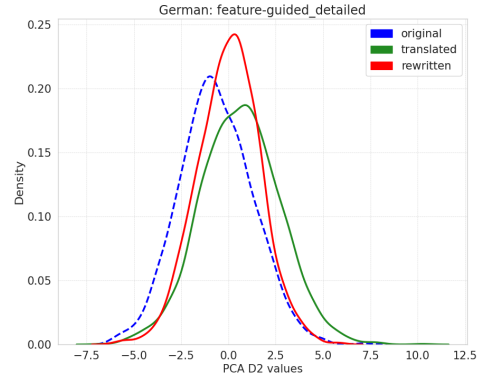


Figure 2: KDE plot for values on the ‘translationese’ dimension of a PCA-transformed data for German translations rewritten in feature-guided detailed mode (on 58 features).

ments that did not exhibit deviations above a 2.5-ratio threshold were not sent to the model. In self-guided modes and re-translation, the model was more willing to dismiss segments as requiring no editing in German than in English. Moreover, for German the number of automatically bypassed segments was close to the number of segments that were skipped in the feature-guided modes, while for English there was a strong contrast in this respect. The translationese filter used in the feature-guided prompt generation considered about 29.28% of HTs into English sufficiently complying with the TL norm, while only less than 1% were not changed in self-guided setups.¹⁴ This means that GPT-4 was more ready to edit a text in English than in German. Note that unchanged segments were included in the data underlying classification results in Table 5 to maintain comparability with the baseline.

Second, we looked into the changes in the feature frequencies in the rewritten segments against the TL non-translations and grouped the features according to their contribution to the task. The *expected* outcome is a reduction of significant deviations from the TL norm. Other possible developments include *no change* compared to the input and some *new trends* absent in HTs. Table 6 has the number of features in each group counted from the full results of statistical tests given in Appendix B.

Table 6 shows that the feature-guided modes had different effectiveness across the translation directions. In German, the expected changes were ob-

¹⁴The full account of these differences can be found in Appendix B, Table 7.

		shift	Feature-guided	
			Min	Detail
DE	expected		3 (0)	6 (3)
	new trend		19 (7)	16 (4)
	no change		36 (8)	36 (8)
EN	expected		15 (2)	16 (3)
	new trend		29 (9)	26 (9)
	no change		14 (4)	16 (3)

Table 6: Analysis of changes in feature frequencies and significance of differences: Number of features by the direction of frequency change after rewriting in feature-guided modes. The number in brackets shows how many of them were among the features addressed in the instructions.

served only for a few features (3 and 6 for Min and Detail modes), while most features remained unaffected (36 for both Min and Detail modes). In English, most features (29 and 26 for Min and Detail modes) demonstrated new deviations from the TL norm. Two-thirds of these emerging trends were over-normalising tendencies, i.e. the features started to deviate from the TL norm in the direction opposite what is typically observed in translations. This effect can hardly be linked to the number of times each feature appeared in the instructions. We hypothesise that the unexpected outcomes were collateral to the other requested transformations which counteracted the specific instructions to favour or avoid specific structures. Except for over-normalisation, the rewritten versions occasionally exhibited deviations on the features where there were no statistical differences between HTs and non-translations.

In almost all cases the non-significant lack or overuse of a specific item was intensified by rewriting. For example, in feature-guided detailed mode on German the number of clauses per sentence (numcls) and specifically of clausal complement without own subjects (xcomp) went further down as compared to HT. In English, the lower frequency of coordinated elements (conj) and higher frequency of simple sentences (simple) reached levels of statistical significance. These deviations, however, were not large and/or consistent enough to build new patterned distinctions between GPT4-edited translations and the TL norm, at least not along the same translationese properties. The rewriting pipeline effectively removed the targeted translationese signals without introducing new deviations, at least those captured by our features. It should be noted that there seems to be a certain limit to the effective number of instructions that

could be passed to the rewriting pipeline. In the reported feature-guided setups, the number of instructions per segment was at most 7 for German and 9 for English, with averages about 2.4 and 2.3, respectively. An attempted alternative approach that generated more instructions per segment was less successful. That approach considered all features with the statistical differences between originals and translations (about 43-44 out of 58 features) if their frequencies for a given translated segment were two standard deviations away from the expected TL norm in the ‘translationese’ direction. This approach generated more varied and longer lists of instructions: the average number of instructions per segment was 3.4, and the number of features addressed in the instruction was twice higher than in the reported approach (21 and 30 for German and English).

Content Preservation. Even if the rewriting pipeline seems to achieve the goals of translationese reduction, we need to make sure that it outputs acceptable translation variants.

		Rewriting Setups				
		RT	Self-guided		Feature-guided	
		-	Min	Detail	Min	Detail
DE	R	0.63	0.87	0.86	0.84	0.85
	QE	0.16	0.49	0.48	0.40	0.44
EN	R	0.85	0.85	0.84	0.80	0.82
	QE	0.46	0.46	0.45	0.33	0.39

Table 7: Average COMET scores for the generated sentences from each of our four rewriting techniques for translationese reduction, compared against the original sentences as references.

Table 7 shows that for German the rewriting setups consistently outperform GPT4-translated sentences in terms of COMET scores for both reference-based (R) and reference-free (QE) evaluations. Specifically, for reference-based (R) evaluation, the COMET scores range from 0.84 to 0.87 across different rewriting setups, indicating a high level of content preservation. This suggests that the rewriting techniques effectively maintain the meaning of the original English sentences. The results for the English pipeline evaluation indicate that (i) GPT-4 is probably much more skilled in translating into English than into German, and that (ii) the rewriting setups, especially in the feature-guided modes, generate less semantically similar translation candidates, even if they seem to be less deviating from the TL norm on some frequency-

based features.

Manual Analysis. The manual analysis by a translation expert was carried out to assess the quality of the re-written output in addition to automatic COMET scores. The human evaluation (Table 8) returned consistently high scores for both accuracy and fluency, giving better results in the German-to-English direction than English-to-German.

		Rewriting Setups				
		RT	Self-guided		Features	
		-	Min	Detail	Min	Detail
DE	A	5.9	5.8	5.8	5.1	5.4
	F	5.4	5.7	5.6	5.4	5.4
EN	A	5.9	5.7	5.9	5.2	5.4
	F	6	6	5.9	5.6	5.8

Table 8: Results of human evaluation for accuracy (A) and fluency (F) in a 1-6 Likert scale.

Both self-guided modes were rated higher than the feature-guided modes. This is in line with the automatic results on content preservation (cf. Table 7). Although the feature-guided instructions were generally followed (92-96% of observations in DE, 96% of observations in EN), it was noticed that they were applied excessively leading to overtransformed renditions as in Example 2 (Appendix D). Human and machine translation preserved one long sentence, showing traces of translationese. The GPT4-rewritten output in the self-guided modes returned 2-3 short sentences whereas the instruction to *make the sentences shorter* resulted in 4 and even 5 shorter sentences for the same input. A similar tendency can be observed in Example 1, where the instruction to *use more adverbial modifiers* in rewriting translations into German in the feature-guided modes resulted in the overuse of adverbials (underlined in the example) and also intensification of the message and therefore decline in accuracy.

6 Conclusion

In this paper, we explore the potential of using LLM-prompts to reduce translationese-related differences between translated and non-translated texts. We evaluate four types of prompts based on either a high-level explanation of the translationese mitigation task or on a micro-managing approach to prompting where the model received segment-tailored instructions to increase or reduce

the frequency of prominent translationese predictors. Our findings demonstrate that GPT-4 **was able to edit human translations to make them less distinguishable** in an automatic classification setup from non-translations in both self-guided and feature-guided LLM-rewriting modes. The best results were seen for English on the prompts containing feature-guided instructions with a linguistic description of special terminology, showing that **the prompting approach benefited from including linguistic knowledge**.

For German the results were less straightforward but the advantages of detailed task information and specific linguistic instructions were visible. The inferior results on the re-translation task provide further evidence in favour of linguistic features for the translationese reduction task. In our experiments, prompting was more effective in the German-to-English translation direction even though the difference between translated and non-translated documents in German was more detectable to start with (as indicated by 5% higher classification results). We can tentatively explain this result by **the language of instruction (English)**, which might prime the model for better performance when generating English output. Future work may need to extend this research by including tasks with instructions in German, especially when the model rewrites translations into German.

Finally, we have seen that even though rewritten translations exhibited some new individual deviations from non-translations on some individual features, they did not coalesce into patterns picked by a classifier. This conclusion is supported by high results from content preservation metrics and from the manual analysis for accuracy and fluency of translations. While our translationese classification-based evaluation shows that LLM-rewriting is effective, in our paper we focus on the tip of the iceberg, i.e. the segments from 200 most contrastive documents in our data set. Furthermore, manual evaluation and, to some extent automatic evaluation, show that content preservation under LLM-rewriting needs more attention, and we will focus on this in our future research.

7 Acknowledgments

This research was supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – SFB 1102 Information Density and Linguistic Encoding, Project-ID 232722074.

References

- Artetxe, Mikel, Gorka Labaka, and Eneko Agirre. 2019. An effective approach to unsupervised machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 194–203, Florence, Italy, jul. Association for Computational Linguistics.
- Artetxe, Mikel, Gorka Labaka, and Eneko Agirre. 2020. Translation artifacts in cross-lingual transfer learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7674–7684, Online, November. Association for Computational Linguistics.
- Banarescu, Laura, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. Abstract meaning representation for sembanking. In *Proceedings of the 7th linguistic annotation workshop and interoperability with discourse*, pages 178–186.
- Biber, Douglas. 1988. *Variations Across Speech and Writing*. Cambridge University Press.
- Borah, Angana, Daria Pylypenko, Cristina Espana-Bonet, and Josef van Genabith. 2023. Measuring spurious correlation in classification: 'Clever Hans' in translationese.
- Chen, Pinzhen, Zhicheng Guo, Barry Haddow, and Kenneth Heafield. 2023. Iterative translation refinement with large language models. *arXiv preprint arXiv:2306.03856*.
- Clark, Jonathan H, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. TyDi QA: a benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454–470.
- Dutta Chowdhury, Koel, Rricha Jalota, Cristina España-Bonet, and Josef Genabith. 2022. Towards debiasing translation artifacts. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3983–3991, Seattle, United States, jul. Association for Computational Linguistics.
- Edunov, Sergey, Myle Ott, Marc Aurelio Ranzato, and Michael Auli. 2020. On the evaluation of machine translation systems trained with back-translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2836–2846, Online, July. Association for Computational Linguistics.
- Evert, Stefan and Stella Neumann. 2017. The impact of translation direction on characteristics of translated texts: a multivariate analysis for English and German. *Empirical Translation Studies: New Methodological and Theoretical Traditions*, 300:47–80.
- Feng, Yutao, Jipeng Qiang, Yun Li, Yunhao Yuan, and Yi Zhu. 2023. Sentence simplification via large language models. *arXiv preprint arXiv:2302.11957*.
- Graham, Yvette, Barry Haddow, and Philipp Koehn. 2020. Statistical power and translationese in machine translation evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 72–81, Online, November. Association for Computational Linguistics.
- Hendy, Amr, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? A comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.
- Hu, Hai and Sandra Kübler. 2021. Investigating translated Chinese and its variants using machine learning. *Natural Language Engineering*, 27(3):339–372.
- Jalota, Rricha, Koel Chowdhury, Cristina España-Bonet, and Josef van Genabith. 2023. Translating away translationese without parallel data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7086–7100.
- Jing, Yingqi and Haitao Liu. 2015. Mean hierarchical distance augmenting mean dependency distance. In Nivre, Joakim and Eva Hajicova, editors, *Proceedings of the Third International Conference on Dependency Linguistics (Depling 2015)*, pages 161–170, 24–26 August.
- Kocmi, Tom and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. In *24th Annual Conference of the European Association for Machine Translation*, page 193.
- Konig, Ekkehard and Volker Gast. 2007. *Understanding English-German Contrasts*. Erich Schmidt Verlag.
- Kunilovskaya, Maria and Ekaterina Lapshinova-Koltunski. 2020. Lexicogrammatic translationese across two targets and competence levels. In Calzolari, Nicoletta, Frederic Bechet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, and And Others, editors, *Proceedings of the 12th Conference on Language Resources and Evaluation (LREC 2020)*, pages 4102–4112. The European Language Resources Association (ELRA).
- Kunilovskaya, Maria and Gloria Corpas Pastor. 2021. Translationese and register variation in English-to-Russian professional translation. In Wang, Vincent X., Defeng Li, and Lily Lim, editors, *New Frontiers in Translation Studies*, pages 133–180. Springer Nature Singapore Pte Ltd.

- Lu, Qingyu, Baopu Qiu, Liang Ding, Liping Xie, and Dacheng Tao. 2023. Error analysis prompting enables human-like translation evaluation in large language models: a case study on ChatGPT. *arXiv preprint arXiv:2303.13809*.
- Nini, Andrea. 2015. Multidimensional Analysis Tagger (v. 1.3).
- Qi, Peng, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: a Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.
- Raunak, Vikas, Arul Menezes, Matt Post, and Hany Hassan. 2023. Do GPTs produce less literal translations? In *The 61st Annual Meeting Of The Association For Computational Linguistics*.
- Ravfogel, Shauli, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. Null it out: guarding protected attributes by iterative nullspace projection. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7237–7256.
- Redelinghuys, Karien. 2016. Levelling-out and register variation in the translations of experienced and inexperienced translators: a corpus-based study. *Stellenbosch Papers in Linguistics*, 45(0):189–220.
- Rei, Ricardo, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585.
- Reif, Emily, Daphne Ippolito, Ann Yuan, Andy Coenen, Chris Callison-Burch, and Jason Wei. 2022. A recipe for arbitrary text style transfer with large language models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 837–848.
- Riley, Parker, Isaac Caswell, Markus Freitag, and David Grangier. 2020. Translationese as a language in “multilingual” NMT. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7737–7746, Online, July. Association for Computational Linguistics.
- Singh, Jasdeep, Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2019. XLDA: cross-lingual data augmentation for natural language inference and question answering. *arXiv preprint arXiv:1905.11471*.
- Suzgun, Mirac, Luke Melas-Kyriazi, and Dan Jurafsky. 2022. Prompt-and-Rerank: a method for zero-shot and few-shot arbitrary textual style transfer with small language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2195–2222.
- Toral, Antonio, Sheila Castilho, Ke Hu, and Andy Way. 2018. Attaining the unattainable? reassessing claims of human parity in neural machine translation. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 113–123, Brussels, Belgium, October. Association for Computational Linguistics.
- Varga, Dániel, Péter Halácsy, András Kornai, Viktor Nagy, László Németh, and Viktor Trón. 2005. Parallel corpora for medium density languages. *International Conference Recent Advances in Natural Language Processing, RANLP*, pages 590–596.
- Vilar, David, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting PaLM for translation: assessing strategies and performance. In *The 61st Annual Meeting Of The Association For Computational Linguistics*.
- Volansky, Vered, Noam Ordan, and Shuly Wintner. 2015. On the features of translationese. *Digital Scholarship in the Humanities*, 30(1):98–118.
- Wein, Shira and Nathan Schneider. 2024. Lost in translationese? Reducing translation effect using abstract meaning representation. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics.
- Yu, Sicheng, Qianru Sun, Hao Zhang, and Jing Jiang. 2022. Translate-train embracing translationese artifacts. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 362–370, Dublin, Ireland, May. Association for Computational Linguistics.
- Zhang, Mike and Antonio Toral. 2019. The effect of translationese in machine translation test sets. In Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 73–81, Florence, Italy, aug. Association for Computational Linguistics.

Appendix A. Linguistic features

Table 7: Types of linguistic information by language level captured with the UD features. The 15 features identified as strong translationese predictors at sentence level for German as a target language appear in bold, for English – in *italics*.

	type	number	list of features [shorthand code]
1	word forms	5	finite verb [fin] , <i>past tense, including conjunctive forms [pastv]</i> , infinitive [inf], passive voice form [aux:pass], deverbal noun [de-verb]
2	word classes	9	noun [nn], <i>personal [ppron]</i> , possessive [poss] , reflexive [self] and demonstrative [demdet] pronouns, adverbial quantifier [advqua], coordinate and subordinate conjunctions ([cconj], [sconj]), adposition [prep]
3	discourse markers	5	<i>adversative [advers]</i> , additive [addit] , causative-consecutive [caus] , temporal-sequential [tempseq] connectives and epistemic stance markers [epist]
4	types of clauses	7	clause with modal predicates [mpred], adjectival clause, including relative clauses [acl], adverbial clause [advcl] , clausal complement with or without own subjects ([ccomp], [xcomp], respectively), asyndetically joined elements in a sentence [paratax] , <i>negative clause [negs]</i>
5	other dependencies	17	adjective in attributive function [amod], adverbial modifier [advmod] , auxiliary verb [aux], appositional modifier [appos], <i>conjunctive relation [conj]</i> , copula verb [cop], three types of relations within multi-word expressions ([compound], [fixed], [flat]), discourse element [discourse], subordinate clause marker [mark], nominal subject [nsubj], direct object [obj], indirect object [iobj] , <i>non-core argument [obl]</i> , numeric modifier [nummod], nominal dependent of a noun [nmod]
6	sentence complexity and word order	10	mean hierarchical distance [mdd] and mean dependency distance [mhd] , <i>number of clauses per sentence [numcls]</i> , ratio of nouns or proper names as core verb arguments to the total of these arguments [nnargs] , ratio of head-verb preceding noun-object to all objects in a clause [vo_noun], inversion in main clause (in affirmative sentences) [vs_noun], ratio of oblique object preceding direct object to clauses with both dependencies [obl_obj], adverbial modifier preceding head-verb to all adverbial modifiers in a clause [adv_verb], any dependencies except subject preceding the main verb [vorfield], prepositional phrases at the end of the finite clauses [nachfield]
7	textual properties	5	lexical type-to-token ratio [ttr] and lexical density [dens] (based on disambiguated content types), number of simple sentences [simple], sentence length [sent_len] and word length [wrlen]
	TOTAL	58	

Appendix B. Changes in feature frequencies and feature importance

Table 7: The expected TL thresholds (i.e. the average feature values in TL originals) and the significance of differences between originals, on the one hand, and HT/rewritten outputs for each feature, on the other hand. The upward and downward departures from the expected TL norm are shown by arrows. The asterisks indicate a lack of statistical significance for the difference based on the two-tailed Mann-Whitney test for unpaired samples. The 15 features identified as strong translationese predictors at sentence level for German as a target language appear in bold, for English – in highlighted rows.

	TL		Rewriting Setups					TL		Rewriting Setups				
	TL	HT	RT	Self	Feature	Min	Det	Min	Det	RT	Self	Feature	Min	Det
	English-to-German							German-to-English						
addit	0.02	↓	↓	↓	↓	↓	↓	0.002	↑	↑	↑	↑	↑*	↓
advcl	0.312	↑	↑	↑	↑	↓	↓	0.552	↑*	↓	↓	↓	↓	↓
advmod	3.327	↓	↓	↓	↓	↓	↓	1.112	↑	↑	↑	↑	↓	↓
caus	0.012	↓	↓	↓	↓	↓	↓	0.002	↑	↑	↑	↑	↑*	↑*
fin	2.673	↓	↓	↓	↓	↓	↓	2.289	↑	↓	↓	↓	↓	↓
iobj	0.153	↑	↑	↑	↑	↓	↓*	0.01	↑	↓*	↓*	↑*	↓*	↑*
mdd	3.512	↑	↓	↓	↓	↓	↓	2.668	↑	↓*	↑*	↓	↓	↓
sent_len	29.222	↓*	↓	↓	↓	↓	↓	27.503	↑	↓	↓	↓	↓	↓
mhd	3.552	↑	↑	↑	↓*	↓	↓	3.857	↓*	↓	↓	↓	↓	↓
nmod	1.257	↑	↑	↑	↑	↓	↓*	1.562	↓	↓	↓	↓	↓	↓
nnargs	0.378	↑	↑	↑	↑	↑	↑	0.584	↓	↓	↓	↓	↓	↓
paratax	0.173	↓	↓	↓	↓	↓	↓	0.059	↑	↑	↑	↑	↓*	↓*
pastv	0.238	↑	↑	↑	↑	↑	↑	0.966	↑	↓*	↓*	↓	↓	↓
poss	0.006	↑	↑	↑	↑	↓	↑*	0.012	↓*	↓	↓	↓*	↓*	↑
ttr	0.958	↑	↑	↑	↑	↑	↑	0.964	↑*	↑	↑	↑	↑	↑
acl	0.407	↑	↑	↑	↓*	↓*	↓*	0.372	↓	↓	↓	↓	↓	↓
advers	0.003	↑*	↑	↑	↑*	↑	↑	0.002	↑	↑	↑	↑	↑	↓*
adv_verb	0.157	↓	↓	↓	↓	↓	↓	0.117	↑	↓	↓	↓	↓*	↓*
advqua	0.023	↓	↓	↓	↓	↓	↓	0.008	↑	↑	↑	↑	↓	↓*
amod	1.288	↑	↑	↑	↓	↓*	↓*	1.702	↑	↓*	↓*	↓	↓	↓
appos	0.163	↓	↓	↓	↓	↓	↓	0.06	↑	↑	↑	↑	↓	↑
aux	0.959	↓	↓	↓	↓	↓	↓	0.853	↑	↓*	↓*	↓	↓	↓
aux:pass	0.24	↑	↑	↑	↓*	↑	↑	0.248	↑	↑	↑	↓*	↓*	↓*
ccomp	0.468	↓	↓	↓	↓	↓	↓	0.294	↑	↑	↑	↓	↓	↓
cconj	0.034	↓	↓	↓	↓	↓	↓	0.035	↓	↓	↓	↓	↓	↓
compoun	0.082	↑*	↓*	↓*	↓*	↓*	↓*	1.012	↓	↓	↓	↓	↓	↓
conj	1.169	↓	↓	↓	↓	↓	↓	1.139	↓*	↓	↓	↓	↓	↓
cop	0.454	↓	↓	↓	↓	↓	↓	0.529	↑	↓	↓*	↓	↓	↓
demdets	0.012	↑	↑	↑	↑	↑	↑	0.017	↓*	↓*	↓*	↓	↓*	↓*
dens	0.41	↓	↓	↓	↓	↑*	↓*	0.423	↓	↓*	↓	↓*	↑	↑*
deverb	0.016	↑	↑	↑	↑	↑	↑	0.025	↓	↓*	↓*	↓	↑	↑
discourse	0.0	↓*	↓*	↓*	↓*	↓*	↓*	0.003	↑	↑*	↑	↑	↓*	↑
epist	0.005	↓	↓	↓	↓	↓	↓	0.003	↑	↑	↑	↑	↓*	↑
fixed	0.011	↓*	↓*	↓*	↓*	↓*	↓*	0.098	↑	↓*	↑	↓*	↓	↓
flat	0.097	↑	↑	↑	↑	↓	↑	0.076	↑	↑	↑	↑	↑	↑
inf	0.008	↓*	↑	↑	↑	↑	↑	0.019	↓	↓*	↑*	↑	↓	↓*
mark	1.03	↓*	↓*	↓*	↓*	↓	↓	1.32	↑	↓*	↓*	↓	↓	↓
mpred	0.6	↓	↓	↓	↓	↓	↓	0.048	↑	↑	↑	↑*	↓*	↑*

nachfeld	0.362	↓*	↓	↓*	↓	↓	↓	0.095	↑*	↓*	↓*	↓*	↓*	↓*
negs	0.012	↓	↓	↓	↓	↓	↓	0.009	↓	↓	↓	↑	↓*	↓*
nn	0.152	↑	↑	↑	↑	↑	↑	0.199	↓	↓	↓	↓	↓*	↓
nsubj	2.356	↓	↓	↓	↓	↓	↓	1.896	↑	↓*	↓*	↓	↓	↓
numcls	1.406	↓*	↓	↓	↓	↓	↓	1.356	↑*	↓	↓	↓	↓	↓
nummod	0.107	↑*	↓*	↓*	↓*	↓*	↓*	0.238	↓	↓	↓	↓	↓	↓
obj	1.273	↑	↑	↑	↓*	↓	↓	1.306	↑*	↓*	↓*	↓*	↓	↓
obl	1.335	↑	↓	↓	↓	↓	↓	1.304	↑	↓	↓	↓	↓	↓
obl_obj	0.097	↑*	↑	↑*	↓*	↓*	↓*	0.07	↓*	↓*	↓*	↓	↓*	↓*
ppron	0.057	↓	↓	↓	↓	↓	↓	0.046	↑	↑	↑	↑	↓	↓*
prep	0.153	↑	↑	↑	↑	↑	↑	0.108	↓	↓	↓	↓	↓	↓
sconj	0.023	↓*	↓*	↓*	↓	↓	↓	0.024	↑	↑	↑	↑	↓	↓*
self	0.003	↑	↑	↑	↑*	↑	↑	0.0	↑	↑	↑	↓*	↑*	↑*
simple	0.273	↓	↓*	↓*	↑	↑	↑	0.273	↑*	↑	↑	↑	↑	↑
tempseq	0.011	↓	↓	↓	↓	↓	↓	0.004	↑	↑	↑	↑	↓*	↓
vo_noun	0.107	↑*	↑*	↓*	↑	↑	↑	0.629	↑*	↓	↓	↓	↓	↓
vorfeld	0.467	↓*	↓	↓	↓	↓*	↓	0.434	↑	↑	↑	↓*	↓	↓
vs_noun	0.044	↑	↑	↑	↑	↑	↑	0.0	↓*	↓*	↓*	↓*	↓*	↓*
wklen	5.6	↑	↑	↑	↑	↑	↑	4.742	↓	↑	↑	↑	↑	↑
xcomp	0.269	↓*	↓*	↓	↓	↓	↓	0.369	↑	↑	↑	↑	↓*	↓*

Table 7: Percentage of segments that did not undergo changes in the re-writing pipeline because no translationese was detected in them either by the model or by feature analysis.

Rewriting Setups					
	RT	Self-guided		Feature-guided	
	–	Min	Detail	Min	Detail
DE	7.92	5.32	0.78	6.24	
EN	0.05	0.05	0.16	29.28	

Appendix C. Examples of prompts by approach and mode

1. **Self-guided approach:** the model has to decide on itself whether a segment contains translationese or not. The same instruction was passed for each pair of segments.

- **Min mode:**

Your task is to re-write a human translation in a more natural way if necessary.

Here is an original English text: “‘In six short months, the presidency has conspired to undermine the Stability Pact, has shown contempt for the European Union’s policy towards Russia and offended Canada.”“

This is its human translation into German: “‘In sechs kurzen Monaten ist es dem Ratsvorsitz gelungen, den Stabilitätspakt zu unterminieren, die Politik der Europäischen Union gegenüber Russland zu missachten und Kanada zu beleidigen.”“

If this translation can be revised to sound more like a text originally produced in the target language, return a revised version. If this translation sounds natural enough, return the input translation.

Do not add any meta-phrases or quotation marks. Do not copy the original text.

- **Detail mode:**

Your task is to reduce translationese in a human translation by re-writing it in a more natural way where possible.

Translationese refers to any regular linguistic features in the translated texts that make them distinct from texts originally produced in the target language, outside the communicative situation of translation. These features are typically detected by statistical analysis and are explained by the specificity of the translation process. Human translators are known to simplify the source language content and to make it more explicit. Translations can exhibit a tendency to conform to patterns which are typical of the target language, making the output less varied than in comparable non-translations in the target language. The more obvious sign of translationese is interference, which can be defined as over-reliance on the intersection of patterns found in source and target languages. Translationese is manifested in the inflated frequencies of specific linguistic items such as function words (especially connectives and pronouns), unusual frequencies of some parts of speech (especially nouns and adverbs) or grammatical forms (especially forms of verbs), in reduced lexical variety and unexpected lexical sequences, in less natural word order, in longer and more complex sentences as well as lack of target language specific items and structures.

Here is an original English text: “‘In six short months, the presidency has conspired to undermine the Stability Pact, has shown contempt for the European Union’s policy towards Russia and offended Canada.”“

This is its human translation into German: “‘In sechs kurzen Monaten ist es dem Ratsvorsitz gelungen, den Stabilitätspakt zu unterminieren, die Politik der Europäischen Union gegenüber Russland zu missachten und Kanada zu beleidigen.”“

If you can detect any translationese deviations in this translation, revise this translation to make it sound less translated and return the revised version. If no translationese is detected, return the input translation.

Do not add any meta-phrases or quotation marks. Do not copy the original text.

2. **Feature-guided approach:** the model is ‘micro-managed’ in how the translation needs to be adapted, if at all. Each pair of segments gets individual instructions, based on features that were found to strongly deviate from the expected TL norm in this translated segment.

- **Min mode:**

Your task is to re-write a human translation in a more natural way.

Here is an original English text: “‘In six short months, the presidency has conspired to undermine the Stability Pact, has shown contempt for the European Union’s policy towards Russia and offended Canada.’”

This is its human translation into German: “‘In sechs kurzen Monaten ist es dem Ratsvorsitz gelungen, den Stabilitätspakt zu unterminieren, die Politik der Europäischen Union gegenüber Russland zu missachten und Kanada zu beleidigen.’”

Re-write this translation following the instructions:

Use pronouns instead of nouns as verbal arguments where possible.

Avoid constructions with indirect objects.

Do not add any meta-phrases or quotation marks. Do not copy the original text.

- **Detail mode:**

Your task is to reduce translationese in a human translation by re-writing it in a more natural, less translated way.

Translationese refers to any properties of translations that make them statistically distinct from texts originally produced in the target language.

Here is an original English text: “‘In six short months, the presidency has conspired to undermine the Stability Pact, has shown contempt for the European Union’s policy towards Russia and offended Canada.’”

This is its human translation into German: “‘In sechs kurzen Monaten ist es dem Ratsvorsitz gelungen, den Stabilitätspakt zu unterminieren, die Politik der Europäischen Union gegenüber Russland zu missachten und Kanada zu beleidigen.’”

Revise this translation following the instructions which reflect deviations of this segment from the expected target language norm:

Use pronouns instead of nouns or proper names as verbal arguments where possible.

Avoid constructions with indirect objects. An indirect object of a verb is any nominal phrase that is an obligatory argument of the verb but is not its subject or direct object. The prototypical example is the recipient (dem Kind) with verbs of exchange: Die Frau gibt dem Kind einen Apfel.

Do not add any meta-phrases or quotation marks. Do not copy the original text.

3. **Re-translation approach:** the model is tasked with re-translation instead of re-writing/revising the existing translation. We keep the prompt as similar as possible to *Self-guided min* setup, while changing the core task formulation. The aim is to compare revised and re-translated versions to see whether the model makes a difference between them.

Your task is to re-translate a human translation to make it more natural in the target language if necessary.

Here is an original English text: “‘In six short months, the presidency has conspired to undermine the Stability Pact, has shown contempt for the European Union’s policy towards Russia and offended Canada.’”

This is its human translation into German: “‘In sechs kurzen Monaten ist es

dem Ratsvorsitz gelungen, den Stabilitätspakt zu unterminieren, die Politik der Europäischen Union gegenüber Russland zu missachten und Kanada zu beleidigen.““
If this translation can be re-translated to sound more like a text originally produced in the target language, return a re-translated version. If this translation sounds natural enough, return the input translation.
Do not add any meta-phrases or quotation marks. Do not copy the original text.

Appendix D. Examples of source segments, human translation as well as all output modes

Ex.	Type	Segments
1	Source:	Does he regard the committee whose advice he rejects as incompetent, unable to judge what is good for animals or dangerous for human health?
	HT:	Hält er den Ausschuß, dessen Empfehlung er ablehnt, für nicht kompetent, für unfähig zu entscheiden, was für die Tiere gut oder für die menschliche Gesundheit gefährlich ist?
	RT:	Betrachtet er den Ausschuss, dessen Ratschlag er ablehnt, als inkompetent und unfähig zu beurteilen, was gut für Tiere oder gefährlich für die menschliche Gesundheit ist?
	Self-guided Min:	Glaubt er etwa, der Ausschuss, dessen Ratschläge er ignoriert, sei inkompetent und unfähig zu beurteilen, was gut für Tiere ist oder der menschlichen Gesundheit schaden könnte?
	Self-guided Detail:	Hält er den Ausschuss, dessen Ratschläge er ablehnt, für inkompetent und unfähig zu beurteilen, was gut für Tiere ist oder dem menschlichen Wohlbefinden schaden könnte?
	Feature-guided Min:	Betrachtet er den Ausschuss, dessen Rat er <u>entschieden</u> ablehnt, <u>tatsächlich</u> als inkompetent, als <u>gänzlich</u> unfähig zu beurteilen, was <u>ausgerechnet</u> für Tiere gut oder <u>potenziell</u> für die menschliche Gesundheit schädlich ist?
	Feature-guided Detail:	Hält er <u>tatsächlich</u> den Ausschuß, dessen Empfehlung er <u>konsequent</u> ablehnt, für <u>gänzlich</u> inkompetent, für <u>absolut</u> unfähig, um zu beurteilen, was <u>letztendlich</u> für die Tiere gut oder <u>womöglich</u> für die menschliche Gesundheit gefährlich ist?
2	Source:	I am pleased that our own legal services – who say Parliament must be very wary of attempts from representatives of other institutions to limit our powers – have confirmed that there is nothing in the rules to stop us having a parliamentary inquiry, so long as the terms of reference differ from any question before the courts in a Member State.
	HT:	Ich bin froh, dass unsere eigenen juristischen Dienste – die sagen, das Parlament müsse vor den Versuchen von Vertretern anderer Institutionen, unsere Befugnisse zu beschneiden, sehr auf der Hut sein – bestätigt haben, dass es in der Geschäftsordnung keine Bestimmung gibt, die uns von einer parlamentarischen Untersuchung abhalten könnte, solange die Angelegenheit nicht bei den Gerichten eines Mitgliedstaates anhängig ist.
	RT:	Ich freue mich, dass unsere eigenen Rechtsdienste - die warnen, dass das Parlament bei Versuchen von Vertretern anderer Institutionen, unsere Befugnisse einzuschränken, äußerst vorsichtig sein muss - bestätigt haben, dass es in den Regeln nichts gibt, was uns daran hindern könnte, eine parlamentarische Untersuchung durchzuführen, solange die Untersuchungsgegenstände sich von denen unterscheiden, die vor Gericht in einem Mitgliedstaat verhandelt werden.
	Self-guided min:	Es freut mich, dass unsere eigenen Rechtsdienste bestätigt haben, dass es keine Regel gibt, die uns daran hindert, eine parlamentarische Untersuchung durchzuführen. Sie warnen jedoch, dass wir sehr vorsichtig sein müssen, wenn Vertreter anderer Institutionen versuchen, unsere Macht zu begrenzen. Dies gilt natürlich nur, solange die Untersuchungsthemen nicht bereits vor den Gerichten eines Mitgliedslandes verhandelt werden.
	Self-guided Detail:	Es freut mich, dass unsere eigenen Rechtsdienste bestätigt haben, dass nichts in der Geschäftsordnung uns daran hindert, eine parlamentarische Untersuchung durchzuführen - solange die Untersuchungsgegenstände nicht bereits vor Gericht in einem Mitgliedsstaat verhandelt werden. Sie warnen das Parlament jedoch, wir sollten wachsam gegenüber Versuchen von Vertretern anderer Institutionen sein, unsere Macht einzuschränken.

Feature-guided Min: Die Freude ist groß, die eigenen juristischen Dienste haben bestätigt. Sie warnen das Parlament, behutsam zu sein. Vertreter anderer Institutionen könnten versuchen, Befugnisse zu beschneiden. Keine Bestimmung in der Geschäftsordnung könnte eine parlamentarische Untersuchung verhindern. Dies gilt, solange die Angelegenheit nicht bei Gericht eines Mitgliedstaates vorliegt.

Feature-guided Detail: Es freut mich, dass ein Bestätigung von den eigenen juristischen Diensten vorliegt. Sie warnen das Parlament, wachsam gegenüber Versuchen von Vertretern anderer Institutionen zu sein, unsere Befugnisse zu beschränken. In der Geschäftsordnung gibt es keine Regel, die eine parlamentarische Untersuchung verhindern könnte. Dies gilt solange der Sachverhalt nicht vor den Gerichten eines Mitgliedstaates liegt.

Translate your Own: a Post-Editing Experiment in the NLP domain

Rachel Bawden¹ Ziqian Peng² Maud Bénard³ Éric de la Clergerie¹
Raphaël Esamotunu³ Mathilde Huguin⁴ Natalie Kübler³ Alexandra Mestivier³
Mona Michelot³ Laurent Romary¹ Lichao Zhu³ François Yvon²

¹Inria, Paris, France

²ISIR, CNRS et Sorbonne Université, Paris, France

³CLILLAC-ARP, Université Paris Cité, Paris, France

⁴INIST, CNRS, Nancy, France

Abstract

The improvements in neural machine translation make translation and post-editing pipelines ever more effective for a wider range of applications. In this paper, we evaluate the effectiveness of such a pipeline for the translation of scientific documents (limited here to article abstracts). Using a dedicated interface, we collect, then analyse the post-edits of approximately 350 abstracts (English→French) in the Natural Language Processing domain for two groups of post-editors: domain experts (academics encouraged to post-edit their own articles) on the one hand and trained translators on the other. Our results confirm that such pipelines can be effective, at least for high-resource language pairs. They also highlight the difference in the post-editing strategy of the two subgroups. Finally, they suggest that working on term translation is the most pressing issue to improve fully automatic translations, but that in a post-editing setup, other error types can be equally annoying for post-editors.

1 Introduction

In most, if not all scientific domains, academic communication and publication activities take place mostly in English (Gordin, 2015). While sharing a common language can be viewed as a facilitating factor in many cases, it also generates tensions, frictions and inequalities (Amano et al.,

2023), and hinders the exposure of science that is not discussed in English. Furthermore, in non-English-speaking countries, it creates a linguistic barrier between the scientific community and the general public that can only amplify misunderstandings and doubts. These issues have motivated calls for changes as expressed in the “Helsinki initiative”.¹ Among the Natural Language Processing (NLP) community, this has motivated the ACL 60-60 special initiative,² aimed at using automatic tools (speech recognition, machine translation (MT)) and resources (multilingual term lists) to help remove these barriers.

In this paper, we report our attempts to use existing MT technologies to translate English scientific documents in the NLP domain into French. As has been well documented for the biomedical domain in the course of the challenges organised at the Conference on Machine Translation since 2016 (see (Neves et al., 2023) for the latest published edition), academic texts pose specific translation challenges, related notably to term translation and the generation of lexically consistent outputs.

Our main goal in this work is to evaluate the current state-of-the-art in MT for the translation of academic NLP texts with a view to using MT to aid NLP authors in the translation and post-editing of abstracts in non-English languages. We base our evaluation on manually post-edited documents by two populations of post-editors: apprentice and well-trained professional translators on the one hand and NLP experts (academics) who are encouraged to post-edit their own articles on the other. The results of this pilot study will help us design and organise a large-scale experiment that

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹<https://www.helsinki-initiative.org/>

²<https://www.2022.aclweb.org/dispecialinitiative>

will ultimately cover all scientific domains. The main questions we aim to answer are the following: (a) what is the effort needed for academics to post-edit automatic translations of texts in their domain of expertise? (b) can we measure the quality of the resulting translations? (c) can we see a difference between existing translation tools? and (d) what are the residual errors that still hinder the translation of academic publications?

To answer these questions, we designed a post-editing protocol aimed to facilitate the voluntary participation of academics in our domain and collected a set of more than 350 abstracts in the NLP domain, which were post-edited once or several times. A large subset of them are also associated with post-editor feedback on the types and severity of errors present. We analysed them in terms of the post-editing effort, measured using HTER (Snover et al., 2006) and studied them in terms of differences in post-editing patterns. We release the resulting corpus and the code for the post-editing interface for future use.³

2 Related Work

Numerous challenges are faced when developing and adapting NLP models to scientific texts, including how to handle domain-specific terminology (including acronyms), and how to ensure coherence at the document level. In recent years, the development of such tools has been a growing area of interest for NLP researchers, with multiple models being published for the scientific and scholarly domains, e.g. SciBERT (Beltagy et al., 2019), PubmedBERT (Gu et al., 2021), Galactica (Taylor et al., 2022) and ScholarBERT (Hong et al., 2023). Specifically for MT, there have been several initiatives, including the recent IWSLT shared task on translating ACL presentations (Agarwal et al., 2023; Salesky et al., 2023). The project that is closest to our own is the COSMAT project (Lambert et al., 2012), whose aim was to develop a pipeline for integrating the translation of scientific documents into the HAL⁴ archiving platform for English–French translation.

A few corpora are available for scientific document translation, covering different types of publications. The biomedical task at WMT, for instance, has produced parallel test sets for a number of years extracted from article abstracts from

PubMed that are available in several languages (Jimeno Yepes et al., 2017; Neves et al., 2018; Bawden et al., 2019; Bawden et al., 2020; Yeganova et al., 2021; Neves et al., 2022; Neves et al., 2023). The SciPar parallel corpus of scientific texts (Roussis et al., 2022) is composed of master’s and doctoral theses across several domains and in multiple languages. S2ORC (Lo et al., 2020) is also multi-discipline and contains monolingual English articles from Semantic Scholar. In the NLP domain, Mariani et al. (2019) compiled and explored a large-scale comparable corpus of about 65k NLP papers from multiple sources, while Tanguy et al. (2020) focus on French, providing a monolingual corpus from the TALN conferences.

Evaluating MT for scientific documents is challenging, particularly as standard metrics may well underestimate the impact of mistranslating scientific terminology if they are considered equal to other words. This is particularly the case for simple surface-based metrics such as BLEU (Papineni et al., 2002), but is also a currently unknown factor for other automatic metrics such as COMET (Rei et al., 2020). According to the human evaluations of the WMT biomedical shared tasks, e.g. (Neves et al., 2023), term translation was one of the factors most impacting judgments of quality over other factors such as style and naturalness. Another problem is the scarcity of parallel texts that can be used for reference-based evaluation. Moreover, those that exist may not be perfect translations, either because there is no guarantee that two abstracts for the same paper in multiple languages were intended to be perfect translations or because the authors are non-native speakers of at least one of the languages. This therefore motivates alternative approaches to evaluation, including reference-less evaluation (for automatic evaluation, this would refer to quality estimation (Specia et al., 2010)) and human evaluation, through post-editing or error annotation for example.

Post-editing has previously been used as a means of evaluating MT quality, either through the time taken to render a text to an acceptable standard or (largely related) through the number of changes that were made, which is the basis for the HTER metric (“Human-targeted Translation Edit Rate”) (Snover et al., 2006; Dorr et al., 2011). This task-based evaluation strategy is less costly both financially and in terms of effort on the part of translators, and can provide clues as to what types of er-

³<https://github.com/ANR-MaTOS/Resources>

⁴<https://hal.science>

rors are being produced by MT systems. It is also a realistic setting in many cases, including ours, where MT systems can be used to provide an initial translation of a text that the author can then modify. For example, the previously mentioned COSMAT project aimed to integrate such software into the publishing platform to facilitate the production of texts in multiple languages by the authors.

3 Data Collection

We collect a corpus of over 20k English NLP titles and abstracts that we translate automatically into French and of which a selection is then post-edited. Basic statistics on the most common types of publications included are in Table 1. As shown in Figure 1, the corpus contains titles and abstracts from various publication types, the most common being conference papers, journal articles, book sections, preprints, reports and books. Once the initial corpus extracted (Section 3.1), each of the titles and abstracts is automatically translated into French using three MT systems (Section 3.3), and finally, we collect post-edits of the translations by translators and members of the NLP community (Section 3.4). The research protocol received a positive evaluation from our university institutional review board. All code and the resulting corpora will be made publicly available. The abstracts belong to the metadata of the articles and therefore can be freely distributed.⁵

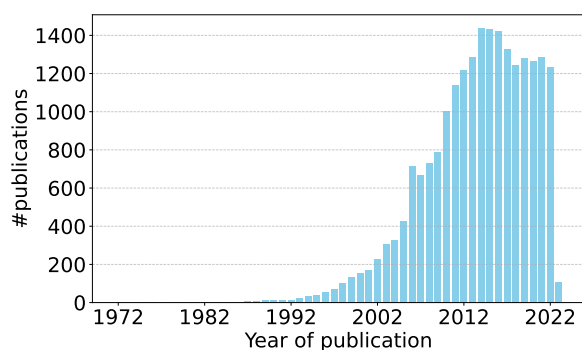


Figure 1: Distribution of publications in the corpus by year.

⁵The metadata of articles published on the HAL platform is under a CC0 licence. It is specified that “HAL’s metadata can be consulted in whole or in part by harvesting in compliance with the intellectual property code.”, pursuant to the so-called French Law for a Digital Republic [Loi n°2016-1321 du 7 octobre 2016 (art. 30)] see, <https://doc.archives-ouvertes.fr/en/legal-aspects/>

Publication by type	#	Avg. #toks	
		title	abstract
Total	21,748	9.86	148.81
Conference paper	14,312	9.70	138.16
Journal article	4,362	10.54	178.68
Book section	1,047	9.38	138.91
Preprint	585	9.87	164.23
Report	506	8.80	150.29
Book	271	9.98	152.15
...			

Table 1: Statistics of the initial NLP corpus overall and for the six most common publication types. Tokens here are defined simply as white-spaced delimited sequences of characters.

3.1 Extracting Scientific Abstracts

Our source texts are English titles and abstracts from scientific publications in the NLP domain from the HAL open archive,⁶ extracted using the dedicated API. In order to select a maximum number of publications with as few non-NLP publications as possible, we carried out the following steps to extract and filter the data: (i) download data from several domains, included the wide domain of “informatics”, (ii) filter to retain only NLP publications, (iii) further filter to remove abstracts that already have a French translation.

Downloading the Data We downloaded the metadata (of which the abstract is one type of information) corresponding to all publications associated with the “computational linguistics” (cs.CL) but also the wider “informatics” domains.

Retaining Only NLP Publications We filter the publications, only keeping those that (i) contain a known keyword in their title, abstract or keyword list or (ii) are published at a known NLP venue.⁷

We check each publication for NLP-specific keywords (in the list of keywords, the title or the abstract). The list was created by taking the set of user-entered keywords for all publications associated with the cs.CL domain, manually filtering it to remove words that could also be relevant to other domains and adding any missing terms based on domain knowledge. This process required manually verifying publications matched with different keywords and removing those that matched with non-NLP publications.

We identify NLP venues by taking the list of conferences, workshops and journals from the

⁶<https://hal.science>

⁷Both the keyword list and venue list can be found at anonymised-link.

ACL anthology corpus (Rohatgi, 2022), adding other known venues and augmenting the list by automatically generating variants of the names (in order to match the various ways authors enter venues), e.g. *13th Nordic Conference of Computational Linguistics (NODALIDA 2001)* also results in *13th Nordic Conference of Computational Linguistics* and *NODALIDA*. We match publications based on the presence of one of the identified venues somewhere in their venue names.

Further Filtering Since the aim is to translate the abstracts into French, we target abstracts that were not originally written in French by filtering out those for which a French abstract exists. This follows the approximation that the presence of a French abstract is likely to indicate that the original language was French.

3.2 Available Metadata

For each article, we collect the following information: title, abstract, list of authors, publication type, venue, date of publication, keywords, language of the text, URL to the paper, licence and the reason for the publication being accepted (out of the filters described above).

3.3 Automatic Translation

We translated the titles and abstracts into French using three commercial neural MT systems: DeepL (professional edition, version 7.5),⁸ Systran Translate (professional edition)⁹ and e-translation (version 12.3).¹⁰ In practice, we concatenated all titles and abstracts into a single file to be translated, separating each article with a token indicating the ID number of the article. We then retrieved the individual translations. Research in contextual MT has shown that when trained properly (this is the case of commercial systems), models have no issue translating multiple sentences at once, especially for short documents such as abstracts (Maruf et al., 2019; Fernandes et al., 2023).

3.4 Manual Post-editing

We developed an online interface to collect post-edits and to provide feedback on MT quality. Users created an account, filling in basic information that

could be useful for future research. They then selected articles to post-edit via the interface and finally gave feedback about the experience. We collected post-edits from two types of post-editors: (i) translators and students¹¹ in translation studies, and (ii) members of the NLP community, who were encouraged, although not forced, to post-edit their own articles.

Post-editor Metadata A condition for participating in the post-editing experiment was fluency in French. Post-editors remained anonymous, but we collected information about their profile that is important for future research, namely their native language(s), other language(s) spoken, the number of years of experience in NLP (<3, 3-10 or 10+) and whether they have previously written an abstract in English and written an abstract in French. We also ask for their general appreciation of MT tools by asking (i) whether they have previously used MT tools to help write scientific articles and (ii) whether they would consider it useful to integrate MT for abstracts into HAL. They can also leave free comments if they wish. Any other information is not available due to anonymity reasons.

Post-editing via the Interface Given that NLP community members were encouraged to post-edit their own publications as experts in the content to be translated, we made sure that they could search the database of publications by ID, keyword (in the title or abstract) and by author name. Otherwise, they could also choose a random publication. To ensure that the same publications were not post-edited too often, publications were presented in a random order in the interface, with a random seed dependent on the ID of the user. Each title and abstract could be post-edited a maximum of three times (once for each MT system). A screenshot of the interface is displayed in Figure 2.

An automatic translation was randomly selected out of the three (the post-editor is unaware of which MT was used). Guidelines were provided on the post-editing page: to modify the text (title and abstract) so that it is clear, understandable and acceptable, as they would do for a journal article written in French. The post-editors could then edit the MT output without a time limit and provide basic feedback on its quality (Figure 3). We also log the time taken to finish post-editing.

¹¹The students worked under the close supervision of their teachers.

⁸<https://deepl.com>

⁹<https://www.systran.net/en/translate>

¹⁰<https://webgate.ec.europa.eu/etranslation>.

Choisir un article à post-éditer

Sélectionnez un article du tableau ci-dessous et cliquez sur ✎ pour faire une nouvelle post-édition. Affinez le choix en cherchant un mot clé, nom d'auteur, année ou identifiant HAL afin de privilégier vos propres articles ou les articles sur certains thèmes : ↕ ✖

Vous pouvez choisir le même article plusieurs fois - une traduction différente sera proposée. Le nombre de post-éditions que vous avez effectuées pour un article donné est indiqué dans la colonne 👤. Le nombre total de post-éditions, tout utilisateur confondu, est dans la colonne 🧑‍🤝‍🧑.

Vous pouvez aussi [choisir un article au hasard !](#)



✎	HAL id	Titre	Auteurs	Année	Lieu	👤	🧑‍🤝‍🧑
✎	1615297	Logic, Formal Linguistics and Computing in France: From Non-reception to Progressive Convergence	Pierre Mounier-Kuhn	2015	3rd International Conference on History and Philosophy of Computing (HaPoC)	0	0
✎	3537323	High-resolution speaker counting in reverberant rooms using CRNN with Ambisonics features	Pierre-Amaury Grumiaux, Srdan Kitic, Laurent Girin, Alexandre Guerin	2021	EUSIPCO 2020 - 28th European Signal Processing Conference (EUSIPCO)	0	0
✎	1557583	Human-Computer Interaction	Peter Forbrig, Fabio Paternò, Annelise Pejtersen	2010	IFIP Advances in Information and Communication Technology	0	0
✎	2880590	Investigating the Impact of Pre-trained Word Embeddings on Memorization in Neural Networks	Aleena Thomas, David Adelani, Ali Davody, Aditya Mogadala, Dietrich Klakow	2020	23rd International Conference on Text, Speech and Dialogue	0	0

Figure 2: Interface for article selection. The instructions read “Select an article from the table below and click on [the pen emoji] to start a new post-edit. Filter your selection by searching for a keyword, author name, year or HAL ID in order to prioritise your own articles or articles from certain themes. [...] You can choose the same article several times - a different translation will be given. The number of times a given article has been post-edited by you is indicated in the column [with the person emoji]. The total number of times it has been post-edited by all users is given in the column [with the people emoji].”

Post-editor Feedback The type of feedback differs depending on the profile of the post-editor. For members of the NLP community, they indicate a single feedback score corresponding to the question “What importance do you give to the MT problems seen?” (as shown in Figure 3), with possible responses “No problem”, “Not very serious (spelling, punctuation, etc.)”, “moderately serious (not interfering with comprehension but not linguistically or stylistically acceptable)” and “serious (interfering with understanding, not faithful to the source)”. As NLP experts are not specialists in manual error annotation, they were presented with four easy-to-use categories. For translators, the question is more detailed, asking for

each error type (faithfulness, grammar, terminology, spelling and punctuation, style, document coherence) whether the problems seen correspond to the same four degrees of quality (“No problem”, “not serious“, moderately serious” or “serious”). In both cases, post-editors can leave a free form comment. The error categories were defined based on the MeLLANGE error typology (Kübler, 2008).

4 NLP Post-edit Corpus

In Table 2, we report basic statistics concerning the post-editing corpus, for documents post-edited by the community (by NLP researchers), for documents post-edited by translator and for two categories combined (all). Given that a single abstract

Post-éditez la traduction d'un titre et d'un résumé dans le domaine du TAL

Instructions :

Modifiez le texte (titre et résumé) pour qu'il soit clair, compréhensible et acceptable, comme vous le feriez pour une publication dans un journal en français (p. ex. la revue TAL). Pour ce faire, merci de ne pas vous servir d'outils de traduction automatique. Dans la mesure du possible, merci de faire cette révision sans vous interrompre pour que la durée enregistrée corresponde au temps effectif de post-édition.

ⓘ Attention : Si vous quittez cette page (en fermant la fenêtre ou en revenant sur la page précédente, vous perdrez les modifications apportées).

Titre :	Investigating alignment interpretability for low-resource NMT
Publié dans :	Machine Translation
Auteurs :	Marcely Zanon Boito, Aline Villavicencio, Laurent Besacier
Année :	2021
ID Hal :	3139744

Résumé d'origine :

Investigating alignment interpretability for low-resource NMT

The attention mechanism in Neural Machine Translation (NMT) models added flexibility to translation systems, and the possibility to visualize soft-alignments between source and target representations. While there is much debate about the relationship between attention and the yielded output for neural models [26, 35, 43, 38], in this paper we propose a different assessment, investigating soft-alignment interpretability in low-resource scenarios. We experimented with different architectures (RNN [5], 2D-CNN [15], and Transformer [39]), comparing them with regards to their ability to produce directly exploitable alignments. For evaluating exploitability, we replicated the Unsupervised Word Segmentation (UWS) task from Godard et al. [22]. There, source words are translated into unsegmented phone sequences. Posterior to training, the resulting soft-alignments are used for producing

Traduction automatique (cliquez pour ouvrir) ▼

Traduction à post-éditer :

Recherche de l'interprétabilité d'alignement pour NMT à faibles ressources

Le mécanisme d'attention dans les modèles de traduction automatique neuronale (NMT) a ajouté de la flexibilité aux systèmes de traduction, et la possibilité de visualiser des alignements souples entre les représentations source et cible. Bien qu'il y ait beaucoup de débat sur la relation entre l'attention et le rendement obtenu pour les modèles neuronaux [26, 35, 43, 38], dans cet article, nous proposons une évaluation différente, en étudiant l'interprétabilité de l'alignement mou dans les scénarios de faibles ressources. Nous avons expérimenté différentes architectures (RNN [5], 2D-CNN [15], et Transformer [39]), en les comparant en ce qui concerne leur capacité à produire des alignements directement exploitables. Pour évaluer l'exploitabilité, nous avons répliqué la tâche de segmentation de mots non supervisés (UWS) de Godard et al. [22]. Là, les mots sources sont traduits

Quelle importance donneriez-vous aux problèmes de traduction constatés ?

- Aucun problème
- Peu grave (orthographe, ponctuation, etc.)
- Moyennement grave (ne gênent pas la compréhension mais linguistiquement ou stylistiquement inacceptables)
- Grave (gênent la compréhension, manquent de fidélité au contenu d'origine)

Remarques libres (optionnel) :

Finaliser

Reinitialiser

[Signaler une erreur technique \(p. ex: pas de résumé\)](#)

Figure 3: Example of the post-editing interface (NLP community member view). The instructions read “Modify the text (title and abstract) so that it is clear, understandable and acceptable, as you would do for a publication in a French journal (e.g. the TAL journal). While post-editing, please do not use machine translation tools. If possible, please complete your post-edition without interruptions so that the registered duration corresponds to the actual time to post-edit... Warning: if you leave this page (by closing the window or going back to the previous page), you will lose your modifications.”. Post-editing is performed without prior sentence segmentation and does not assume that the source and target texts have matching number of sentences.

can be translated multiple times (using different MT systems), we distinguish the statistics concerning the number of abstracts that have been post-edited and the number of translations that have been post-edited (a translation being specific to a particular abstract).

Type	comm.	trans.	all
PEs	95	242	337
Abstracts w/ PEs	91	241	322
Translations w/ PEs	73	240	313
Abstracts w/ several PEs	17	2	55
Translations w/ several PEs	4	1	30

Table 2: Basic statistics concerning the number of post-editions (PEs) by NLP **community** members, by **translators** and by either group (**all**). Among abstracts and translations with several PEs, 46 distinct abstracts are post-edited by both groups, and 28 different translations are post-edited by both.

Concerning the post-editors, there were 4 translators (3 of whom were native French speakers) and 16 NLP experts (13 of whom were native French speakers) and whose experience in NLP ranged from 10+ years (4 users), to 3-10 years (7 users), to under 3 years (5 users).

5 Analysis of the Post-edit Corpus

5.1 Evaluation setup

We primarily base our evaluation of post-editing efforts on the computation of HTER (Human Translation Edit Rate) (Snover et al., 2006), which corresponds to a modified edit distance between the automatic translation and its revised version. We compute HTER with SacreBLEU’s implementation¹² (Post, 2018). Scores are computed separately for each whole abstract (viewed as one long line of text) then broken down by post-editor type and averaged over the corresponding documents.

We also report BLEU score differences between the original and modified abstracts, also computed with SacreBLEU,¹³ in order to judge how much or little the translations had to be edited to be deemed acceptable. These scores rely on corpus-level statistics, again computed on a per-document basis.¹⁴ Measures of post-editing time were also recorded, but we deem them insufficiently reliable

¹²Version 13.5.

¹³We use the default signature for BLEU: nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1.

¹⁴Sometimes referred to as “document-level” BLEU. Note that BLEU scores cannot be used here to compare translation quality across populations, given that they do not use the same set of references.

in our experimental setting to perform any fine-grained analysis.

To evaluate the quality of translation without references, we use Comet-QE (Rei et al., 2020), which relies on distances between continuous space representations of source and target texts.¹⁵

5.2 Results

A first observation is that in our conditions, the automatic translations are mostly of high quality, with an average HTER of 10.7 (BLEU=85.6). Another indication of this high quality is that 13 documents (out of 337) were left entirely unchanged. For the NLP community group, revising an abstract took less than 10 minutes on average.

Comparing the community and translators A more detailed analysis of the post-editing results is illustrated in Figure 4. Two interesting trends can be seen: (a) the distribution of efforts is more concentrated for translators than for the NLP community, (b) the translators also tend to make smaller changes to the translation than the NLP community (HTER=8.0 vs. HTER=18.2), a quite significant difference. This is also obvious when considering the 90% percentile of HTER values (17.2 vs. 32.7). These differences may reveal differences in the way the task was perceived by each population: while translators tend to follow established post-editing guidelines and remain as close as possible to the original MT, field experts are more inclined to rewrite substantial portions of the abstracts.

Without human references, it is difficult to assess the quality of the resulting translations. Computing Comet-QE scores before and after post-editing however reveals a very small improvement (see Table 3). This hints at the lack of sensitivity of QE scores for high-quality translations.

	MT outputs	Post-editions
Translators	76.3	77.0
NLP experts	77.8	78.6

Table 3: Comet-QE(x100) scores of MT outputs and their post-edited versions for each group of post-editors.

Another measure is to take the professional translations as references for the 28 MT outputs post-edited by both groups. For this subset of abstracts, the BLEU score is 76.7 (HTER=18.4).

¹⁵The model is Unbabel/wmt22-cometkiwi-da. See <https://unbabel.github.io/COMET/>.

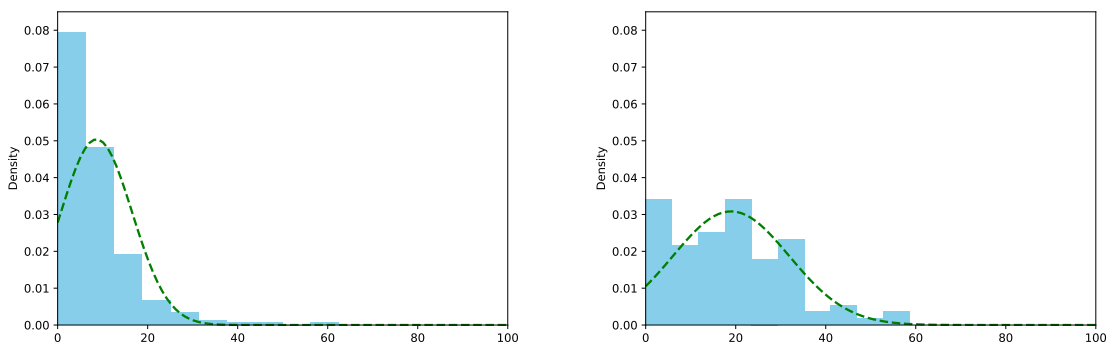


Figure 4: Distribution of HTER scores for translators (left) and NLP experts (right).

Comparison of MT systems We now turn to our third question, which concerns the differences between MT systems. Table 4 reports the average scores for each system and post-editor group and Figure 5 plots the corresponding distributions.

Group	DeepL	eTranslation	Systran
Translators	90.4/7.1 (N=90)	85.9/10.6 (N=79)	87.3/8.5 (N=73)
Experts	81.7/13.8 (N=34)	68.3/24.4 (N=28)	73.1/19.7 (N=33)

Table 4: BLEU/hTER for each post-editor group and system. The number of abstracts for each category is given in brackets.

The scores in Table 4 show clear preferences, with DeepL yielding the smallest post-editing effort, while e-Translation consistently leads to more corrections. These differences are particularly strong for the NLP experts group. These observations are only partly confirmed by a two-sided student T-test for all pairs of systems: out of 6 comparisons, the only significant differences at $p=0.05$ are for DeepL vs. eTranslation for both groups, while Systran cannot be viewed as significantly worse than DeepL, nor significantly better than eTranslation.

Qualitative analysis of errors For the translator group, we analyse the post-editor feedback concerning translation errors for the 7 broad error categories introduced in Section 3.4. We report the corresponding statistics in Table 5.

A first observation is the consistency of these judgements: for each error type, more severe errors tend to yield more edits, with some small inconsistencies (e.g. terminology errors with severity 2 and 3). Looking now at error types, we see that that grammar, style, and punctuation errors are

mostly associated with the lowest level of severity. This is expected given the very high fluidity of MT outputs. The same trend is observed for faithfulness and coherence errors, which tend to get rarer as the severity level increases. Terminology errors exhibit the reverse trend and are mostly associated with the highest level of severity. However, looking now at the post-editing effort, we observe at all severity levels that fixing term errors always yields the lowest HTER scores, while fixing grammar errors almost always yields the highest ones.

Qualitative differences between groups Finally, we carry out a small qualitative analysis of the way the two groups post-edit MT. A few interesting examples are given in Table 6, corresponding to cases where a) one group left the MT output unchanged while the other had high HTER and BLEU scores at the sentence level or b) both groups had high but different HTER and BLEU scores at the sentence level. We note that, while there are a few cases in which the translators corrected an MT error that the community seem to overlook (“traduction automatique de neurone” (literally *machine translation of neurons*) in Example 3), in most cases, the community group seems to produce better post-edited texts than the translators. NLP experts seem to better master specialised terminology (“analyse syntaxique en constituants lexicalisés” in Example 2), specialised phraseology (e.g. “les modèles sont entraînés” *models are trained* instead of “les modèles sont formés”, literally *models are educated*, in Example 4 (source text: “All our models are trained without the need of cross-modal labeled translation data.”)), as well as domain conventions (in Example 1 the acronym “CoMMuTE” is associated with

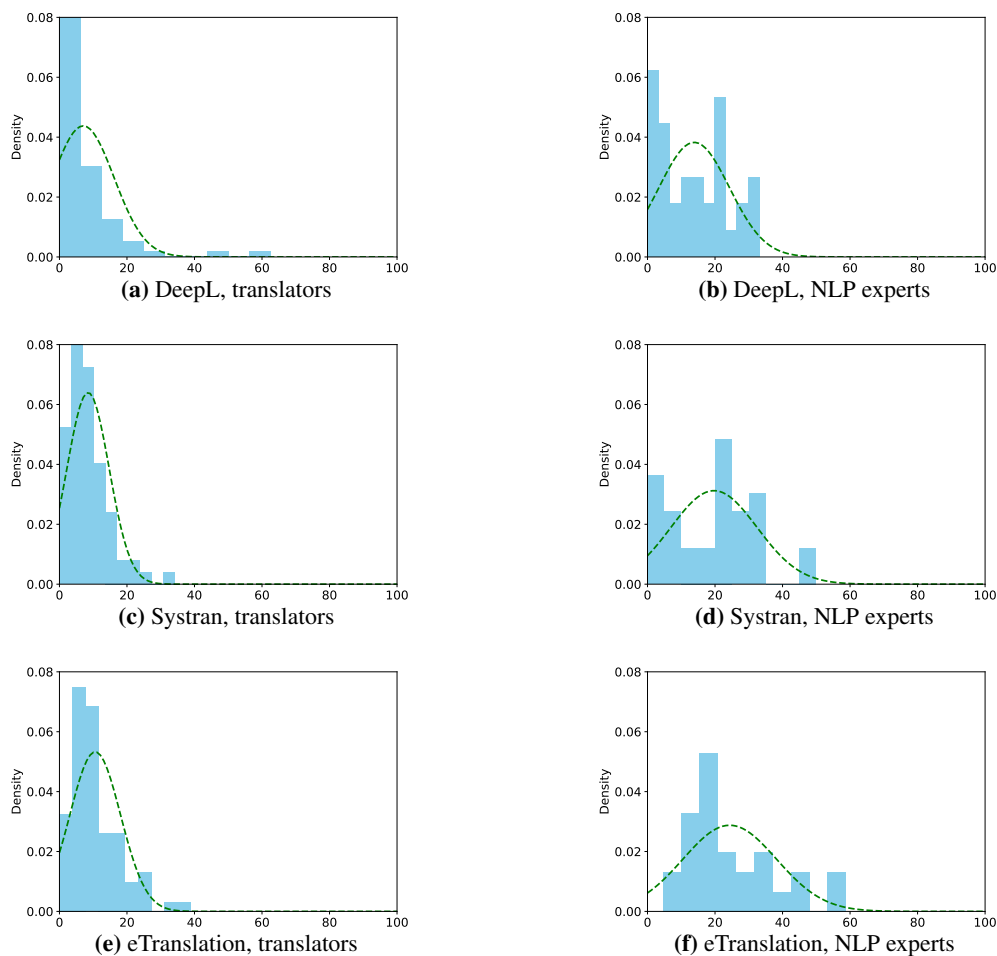


Figure 5: Distribution of HTER scores across systems and post-editor groups.

Problem	Severity level											
	1			2			3			4		
	N	BLEU	hTER	N	BLEU	hTER	N	BLEU	hTER	N	BLEU	hTER
Grammar	148	90.2	7.0	67	84.8	10.7	14	86.4	10.3	12	79.5	16.4
Spelling & Punct.	130	89.5	7.4	60	86.2	10.0	35	87.4	8.8	16	82.8	13.7
Document	127	91.2	6.0	43	84.9	10.7	39	85.2	10.9	33	82.7	13.5
Style	68	91.1	6.1	84	88.7	8.1	57	85.9	10.4	32	82.7	12.7
Faithfulness	123	92.1	5.4	58	84.9	10.4	35	80.8	14.7	25	84.5	12.4
Terminology	34	95.0	3.2	43	88.4	8.3	73	87.5	8.9	91	85.4	10.7

Table 5: Post-edition efforts evaluated according to the number (N, left), BLEU (middle) and hTER (right) of translations associated with different severity levels (from 1 to 4) for translation problems reported by translators in their feedback.

the full term in brackets, which better conforms to the domain conventions than the solution the translator adopted, i.e. translating the full term). Experts also seem to take more freedom in rearranging constituents and rewriting sentences (Example 5), where translators seem to follow the source sentence structure more closely, a behaviour that is also reflected in the automatic metric scores.

6 Conclusion and Future Work

In this paper, we report the results of a pilot study aimed at evaluating the quality of commercial MT systems for scholarly documents (abstracts) in the NLP domain (for English→French). This study explores a realistic scenario, where domain experts post-edit in their mother tongue their own texts (in English, supposedly their L2). We compare against the use of translators with a partial knowledge of the target domain to perform the same task.

MT	NLP experts post-edits	Translator’s post-edits
<i>1- We also release CoMMuTE, a Contrastive Multilingual Multimodal Translation Evaluation dataset, composed of ambiguous sentences and their possible translations, accompanied by disambiguating images corresponding to each translation.</i>		
Nous publions également CoMMuTE, un ensemble de données d’évaluation de la traduction multimodale multilingue contrastive, composé de phrases ambiguës et de leurs traductions possibles, accompagnées d’images désambiguïsantes correspondant à chaque traduction.	Nous publions également le jeu de données CoMMuTE (Contrastive Multilingual Multimodal Translation Evaluation), composé de phrases ambiguës et de leurs traductions possibles, accompagnées d’images visant à leur désambiguïsation et correspondant à chaque traduction.	Nous publions également CoMMuTE, un ensemble de données d’évaluation de la traduction multimodale multilingue contrastive, composé de phrases ambiguës et de leurs traductions possibles, accompagnées d’images désambiguïsantes correspondant à chaque traduction.
<i>2- Multilingual Lexicalized Constituency Parsing with Word-Level Auxiliary Tasks</i>		
Lexicalized Constituency Parsing multilingue avec des tâches auxiliaires de niveau Word	Tâches auxiliaires au niveau des mots pour l’analyse syntaxique en constituants lexicalisés multilingue	Analyse syntaxique de constituants lexicaux multilingues avec tâches auxiliaires au niveau des mots
<i>3- Priming Neural Machine Translation</i>		
Amorçage de la traduction automatique de neurones	Amorçage de la traduction automatique de neurones	Amorçage de la traduction automatique neuronale
<i>4- All our models are trained without the need of cross-modal labeled translation data.</i>		
Tous nos modèles sont formés sans avoir besoin de données de traduction étiquetées intermodales.	Tous nos modèles sont entraînés sans que des données de traduction intermodales annotées soient nécessaires.	Tous nos modèles sont formés sans avoir besoin de données de traduction étiquetées intermodales.
<i>5- On the SPMRL dataset, our parser obtains above state-of-the-art results on constituency parsing without requiring either predicted POS or morphological tags, and outputs labelled dependency trees.</i>		
Sur l’ensemble de données SPMRL, notre analyseur obtient ci-dessus des résultats de pointe sur l’analyse des circonscriptions sans nécessiter une prévision de POS ou d’étiquettes morphologiques, et des sorties marquées d’arbres de dépendance.	Sur l’ensemble de données SPMRL, notre analyseur obtient des résultats supérieurs à l’état de l’art en analyse syntaxique en constituants sans nécessiter de parties du discours prédites ni d’étiquettes morphologiques prédites, et permet de construire des arbres syntaxiques en dépendances étiquetées.	Sur l’ensemble de données SPMRL, notre analyseur obtient des résultats supérieurs à l’état de l’art sur l’analyse des constituants sans nécessiter de prédiction des parties du discours ou des étiquettes morphologiques, ni des sorties marquées d’arbres de dépendance.

Table 6: Comparison of experts’ and translators’ post-edits. Source texts are shown in grey.

Using a dedicated interface adapted for the two populations of post-editors, we collected and analysed approximately 350 abstracts and their post-edited versions. Our main result is that the automatic outputs are already quite satisfactory, as acknowledged by a low average post-editing effort (see also (Sebo and de Lucia, 2024)). We also observed that domain experts tend to deviate more from the original text than translators, the two categories displaying different patterns of post-edits. This study also confirmed the prevalence and severity of terminology errors, while other error types are comparatively rarer or less severe. All resources and analyses will be released to the community.

In the future, we plan to both continue analysis of the data, in particular concerning term use and to reproduce this small-scale experiment with another group of academics from a different scien-

tific background. This will however require finding better ways to incentivise researchers to participate in post-editing activities.

Acknowledgments

This research was funded by the French Agence Nationale de la Recherche (ANR) under the project MaTOS - “ANR-22-CE23-0033-03”. R. Bawden’s participation was also partly funded by her chair in the PRAIRIE institute funded by the ANR as part of the “Investissements d’avenir” programme under the reference ANR-19-P3IA-0001.

References

Agarwal, Milind, Sweta Agrawal, Antonios Anastopoulos, Luisa Bentivogli, Ondřej Bojar, Claudia Borg, Marine Carpuat, Roldano Cattoni, Mauro Cettolo, Mingda Chen, William Chen, Khalid Choukri,

- Alexandra Chronopoulou, Anna Currey, Thierry Declerck, Qianqian Dong, Kevin Duh, Yannick Estève, Marcello Federico, Souhir Gahbiche, Barry Haddow, Benjamin Hsu, Phu Mon Htut, Hirofumi Inaguma, Dávid Javorský, John Judge, Yasumasa Kano, Tom Ko, Rishu Kumar, Pengwei Li, Xutai Ma, Prashant Mathur, Evgeny Matusov, Paul McNamee, John P. McCrae, Kenton Murray, Maria Nadejde, Satoshi Nakamura, Matteo Negri, Ha Nguyen, Jan Niehues, Xing Niu, Atul Kr. Ojha, John E. Ortega, Proyag Pal, Juan Pino, Lonneke van der Plas, Peter Polák, Elijah Rippeth, Elizabeth Salesky, Jiatong Shi, Matthias Sperber, Sebastian Stüker, Katsuhito Sudoh, Yun Tang, Brian Thompson, Kevin Tran, Marco Turchi, Alex Waibel, Mingxuan Wang, Shinji Watanabe, and Rodolfo Zevallos. 2023. Findings of the IWSLT 2023 Evaluation Campaign. In *Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023)*, pages 1–61, Toronto, Canada (in-person and online), July. Association for Computational Linguistics.
- Amano, Tatsuya, Valeria Ramírez-Castañeda, Violeta Berdejo-Espinola, Israel Borokini, Shawan Chowdhury, Marina Golivets, Juan David González-Trujillo, Flavia Montaña-Centellas, Kumar Paudel, Rachel Louise White, and Diogo Veríssimo. 2023. The manifold costs of being a non-native English speaker in science. *PLoS biology*, 21(7):e3002184, July.
- Bawden, Rachel, Kevin Bretonnel Cohen, Cristian Grozea, Antonio Jimeno Yepes, Madeleine Kittner, Martin Krallinger, Nancy Mah, Aurelie Neveol, Mariana Neves, Felipe Soares, Amy Siu, Karin Verspoor, and Maika Vicente Navarro. 2019. Findings of the WMT 2019 biomedical translation shared task: Evaluation for MEDLINE abstracts and biomedical terminologies. In *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 29–53, Florence, Italy, August. Association for Computational Linguistics.
- Bawden, Rachel, Giorgio Maria Di Nunzio, Cristian Grozea, Inigo Jauregi Unanue, Antonio Jimeno Yepes, Nancy Mah, David Martinez, Aurélie Névél, Mariana Neves, Maite Oronoz, Olatz Perez-de Viñaspre, Massimo Piccardi, Roland Roller, Amy Siu, Philippe Thomas, Federica Vezzani, Maika Vicente Navarro, Dina Wiemann, and Lana Yeganova. 2020. Findings of the WMT 2020 biomedical translation shared task: Basque, Italian and Russian as new additional languages. In *Proceedings of the Fifth Conference on Machine Translation*, pages 660–687, Online, November. Association for Computational Linguistics.
- Beltagy, Iz, Kyle Lo, and Arman Cohan. 2019. SciBERT: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3615–3620, Hong Kong, China, November. Association for Computational Linguistics.
- Dorr, Bonnie, Joseph Olive, John McCary, and Caitlin Christianson. 2011. *Machine Translation Evaluation and Optimization*, pages 745–843. Springer New York, New York, NY.
- Fernandes, Patrick, Kayo Yin, Emmy Liu, André Martins, and Graham Neubig. 2023. When does translation require context? a data-driven, multilingual exploration. In Rogers, Anna, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 606–626, Toronto, Canada, July. Association for Computational Linguistics.
- Gordin, Michael D. 2015. *Scientific Babel How Science Was Done Before and After Global English*. University of Chicago Press.
- Gu, Yu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. *ACM Trans. Comput. Healthcare*, 3(1), oct.
- Hong, Zhi, Aswathy Ajith, James Pauloski, Eamon Duede, Kyle Chard, and Ian Foster. 2023. The diminishing returns of masked language models to science. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1270–1283, Toronto, Canada, July. Association for Computational Linguistics.
- Jimeno Yepes, Antonio, Aurélie Névél, Mariana Neves, Karin Verspoor, Ondřej Bojar, Arthur Boyer, Cristian Grozea, Barry Haddow, Madeleine Kittner, Yvonne Lichtblau, Pavel Pecina, Roland Roller, Rudolf Rosa, Amy Siu, Philippe Thomas, and Saskia Trescher. 2017. Findings of the WMT 2017 biomedical translation shared task. In *Proceedings of the Second Conference on Machine Translation*, pages 234–247, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Kübler, Natalie. 2008. A Comparable Learner Translator Corpus: creation and use. In *Proceedings of the 2nd Workshop on Building and Using Comparable Corpora*, pages 73–78, Marrakech, Morocco.
- Lambert, Patrik, Holger Schwenk, and Frédéric Blain. 2012. Automatic translation of scientific documents in the HAL archive. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 3933–3936, Istanbul, Turkey, May. European Language Resources Association (ELRA).
- Lo, Kyle, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. S2ORC: The semantic scholar open research corpus. In *Proceedings of the*

- 58th Annual Meeting of the Association for Computational Linguistics, pages 4969–4983, Online, July. Association for Computational Linguistics.
- Mariani, Joseph, Gil Francopoulo, and Patrick Paroubek. 2019. The NLP4NLP corpus (i): 50 years of publication, collaboration and citation in speech and language processing. *Frontiers in Research Metrics and Analytics*, 3.
- Maruf, Sameen, André F. T. Martins, and Gholamreza Haffari. 2019. Selective attention for context-aware neural machine translation. In Burstein, Jill, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3092–3102, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Neves, Mariana, Antonio Jimeno Yepes, Aurélie Névéol, Cristian Grozea, Amy Siu, Madeleine Kitzner, and Karin Verspoor. 2018. Findings of the WMT 2018 biomedical translation shared task: Evaluation on Medline test sets. In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 324–339, Belgium, Brussels, October. Association for Computational Linguistics.
- Neves, Mariana, Antonio Jimeno Yepes, Amy Siu, Roland Roller, Philippe Thomas, Maika Vicente Navarro, Lana Yeganova, Dina Wiemann, Giorgio Maria Di Nunzio, Federica Vezzani, Christel Gerardin, Rachel Bawden, Darryl Johan Estrada, Salvador Lima-lopez, Eulalia Farre-maduel, Martin Krallinger, Cristian Grozea, and Aurelie Neveol. 2022. Findings of the WMT 2022 biomedical translation shared task: Monolingual clinical case reports. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 694–723, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Neves, Mariana, Antonio Jimeno Yepes, Aurélie Névéol, Rachel Bawden, Giorgio Maria Di Nunzio, Roland Roller, Philippe Thomas, Federica Vezzani, Maika Vicente Navarro, Lana Yeganova, Dina Wiemann, and Cristian Grozea. 2023. Findings of the WMT 2023 biomedical translation shared task: Evaluation of ChatGPT 3.5 as a comparison system. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 43–54, Singapore, December. Association for Computational Linguistics.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October. Association for Computational Linguistics.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.
- Rohatgi, Shaurya. 2022. Acl anthology corpus with full text. Github.
- Roussis, Dimitrios, Vassilis Papavassiliou, Prokopis Prokopidis, Stelios Piperidis, and Vassilis Katsouras. 2022. SciPar: A collection of parallel corpora from scientific abstracts. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 2652–2657, Marseille, France, June. European Language Resources Association.
- Salesky, Elizabeth, Kareem Darwish, Mohamed Al-Badrashiny, Mona Diab, and Jan Niehues. 2023. Evaluating multilingual speech translation under realistic conditions with resegmentation and terminology. In *Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023)*, pages 62–78, Toronto, Canada (in-person and online), July. Association for Computational Linguistics.
- Sebo, Paul and Sylvain de Lucia. 2024. Performance of machine translators in translating French medical research abstracts to English: A comparative study of DeepL, Google Translate, and CUBBITT. *PLOS ONE*, 19(2):1–13, February. Publisher: Public Library of Science.
- Snover, Matthew, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A Study of Translation Edit Rate with Targeted Human Annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.
- Specia, Lucia, Dhvaj Raj, and Marco Turchi. 2010. Machine translation evaluation versus quality estimation. *Machine Translation*, 24(1):39–50, March.
- Tanguy, Ludovic, Cécile Fabre, and Yoann Bard. 2020. Impact de la structure logique des documents sur les modèles distributionnels : expérimentations sur le corpus TALN (impact of document structure on distributional semantics models: a case study on NLP research articles). In *Actes de la 6e conférence conjointe Journées d'Études sur la Parole (JEP, 33e édition), Traitement Automatique des Langues Naturelles (TALN, 27e édition), Rencontre des Étudiants Chercheurs en Informatique pour le Traitement Automatique des Langues (RÉCITAL,*

22e édition). *Volume 2 : Traitement Automatique des Langues Naturelles*, pages 122–135, Nancy, France, 6. ATALA et AFCP.

Taylor, Ross, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science.

Yeganova, Lana, Dina Wiemann, Mariana Neves, Federica Vezzani, Amy Siu, Inigo Jauregi Unanue, Maite Oronoz, Nancy Mah, Aurélie Névéol, David Martinez, Rachel Bawden, Giorgio Maria Di Nunzio, Roland Roller, Philippe Thomas, Cristian Grozea, Olatz Perez-de Viñaspre, Maika Vicente Navarro, and Antonio Jimeno Yepes. 2021. Findings of the WMT 2021 biomedical translation shared task: Summaries of animal experiments as new test set. In *Proceedings of the Sixth Conference on Machine Translation*, pages 664–683, Online, November. Association for Computational Linguistics.

Pre-task perceptions of MT influence quality and productivity: the importance of better translator-computer interactions and implications for training

Vicent Briva-Iglesias

SALIS, D-REAL

Dublin City University

vicent.brivaiglesias2@mail.dcu.ie

Sharon O'Brien

SALIS

Dublin City University

sharon.obrien@dcu.ie

Abstract

This paper presents a user study with 11 professional English-Spanish translators in the legal domain. We analysed whether negative or positive translators' pre-task perceptions of machine translation (MT) being an aid or a threat had any relationship with final translation quality and productivity in a post-editing workflow. Pre-task perceptions of MT were collected in a questionnaire before translators conducted post-editing tasks and were then correlated with translation productivity and translation quality after an Adequacy-Fluency evaluation. Each participant translated 13 texts over two consecutive weeks, accounting for 120,102 words in total. Results show that translators who had higher levels of trust in MT and thought that MT was not a threat to the translation profession reported higher translation quality and productivity. These results have critical implications: improving translator-computer interactions and fostering MT literacy in translation training may be crucial to reducing negative translators' pre-task perceptions, resulting in better translation productivity and quality, especially adequacy.

1 Introduction

MT has become an undisputed element of today's workflows in the language services industry (ELIS Research, 2023). Different studies suggest that improvements in these systems over time have allowed translators to see their productivity increase

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

without a negative impact on the quality produced and, therefore, most research in the field of MT has focused on estimating the productivity and quality of MT systems (Moorkens et al., 2018a; Rossi and Carré, 2022). However, this adoption of MT is not always accompanied by positive user feedback, as some translators have shown little satisfaction in working and interacting with MT through post-editing workflows, either because of a reduction in pay, a sense of dehumanisation of the translation process, or the commodification and uberisation of the language services industry (Moorkens, 2020; Firat, 2021; Cadwell et al., 2018).

In the Translation Studies and the MT fields, research centered on analysing human factors in today's translator-computer interactions is still relatively limited, and the perceptions, user experiences (UX) or feelings of MT users, or even whether these feelings and experiences have any effect on their interactions, have been scarce (Koponen et al., 2020; Karakanta et al., 2022; Briva-Iglesias and O'Brien, 2023; Briva-Iglesias et al., 2023; Guerberof Arenas et al., 2021). In this context, we present the results of a study (part of a larger project) (Briva-Iglesias, 2024) that explores whether translators' pre-task perceptions of MT have any relationship with final translation quality and productivity. Below, we first present work related to our study, then we outline the methodology and, finally, results are described and discussed.

2 Related Work

In the last decades of research in natural language processing (NLP), the focus has been on making technical advancements, mainly by increasing the size of the language models and the computational power used to obtain better results (Brown et al., 2020), but often neglecting the repercus-

sions or risks that this path has or may have on humans (Bender et al., 2021; Shneiderman, 2022). Research in translation technologies has followed a parallel path to NLP research, and most studies have focused on evaluating the quality of MT or comparing different MT paradigms (Drugan, 2013; Moorkens et al., 2018a; Rossi and Carré, 2022). Such research is necessary, but technical changes should also be accompanied by socio-technical studies and their impact on users. Olohan (2011) criticized this path by commenting that “the human and organisational aspects are not addressed at all, or only implicitly, [...] when the system is being developed”.

This has meant that the study of human factors and their interaction with technology has lagged behind and received less attention in translation technology research. However, it has not been completely forgotten. For example, Gaspari et al. (2014) analysed the perceptions of 438 users of online MT systems, and the majority of participants commented that they were not happy with the results, especially with the quality offered. Moorkens et al. (2018b) studied the perception of post-editing effort in the literary field, and collected the data with questionnaires and short interviews, which they then analysed qualitatively. Through a questionnaire completed by 1850 people, O’Brien et al. (2017) investigated how translators interacted with CAT tools, and found that there were certain levels of cognitive friction and that some functionalities of CAT tools irritated them.

Not only freelance or corporate translators have received the attention of academia, but also translators in governmental organisations and international institutions. Rossi and Chevrot (2019) surveyed French translators at the European Commission to analyse the level of acceptance of MT, and suggested that fear of the technology was the element that hindered its adoption. Cadwell, O’Brien, and Teixeira (2018) conducted a similar study, comparing the level of MT uptake of in-house and institutional translators, sharing similar results.

Translation in a migration context has also received attention, as multilingual communication is key in crisis scenarios (Piller et al., 2020), and Pérez-Macías, Ramos, and Rico (2020) analysed the perceptions of MT and post-editing of translators in a migration context, which were negative in general terms.

In contrast, Koponen et al. (2020) focused on

the audiovisual domain and analysed what 12 professional translators thought about MT and what was their UX after post-editing subtitles. The resulting comments ranged from negative to neutral. These results are in line with other research on audiovisual translation, post-editing and UX, where translators do not view post-editing in subtitling projects favourably (Etchegoyhen et al., 2018; Matusev et al., 2019; Karakanta et al., 2022). In a similar vein, Briva-Iglesias, O’Brien, and Cowan (2023) analysed the MTUX of translators in the legal domain to see what translators preferred from two different post-editing modalities.

However, despite having found research on the perceptions that translators have of MT, the aforementioned studies are exclusively descriptive of participant’s perceptions and did not analyse whether these perceptions have any relationship with the quality of the final text or the productivity of translators. This is the gap that this article aims to fill.

In cognitive science, multiple studies show that past experiences and perceptions have a great impact and are a determinant for future beliefs, attitudes and behaviours (Albarracín, 2021; Albarracín and Wyer, 2000). In Translation Studies, de Almeida (2013) suggested that positive perceptions towards MT had an impact on post-editing effort, and Stasimioti and Sosoni (2019) reported that training in MT changed perceptions of MT and post-editing.

Hence, if translators are not happy with their past interactions with MT, and if, before starting a post-editing assignment, they already have a negative opinion about that future interaction (pre-task perceptions), what will the consequences be for the final product (that is, the translation)? Are we in a vicious circle in which translators’ negative pre-task perceptions of MT affect the final quality of the translation and/or their productivity?

3 Methodology

The research question we address in this paper is: *Do translators’ (positive or negative) pre-task perceptions of MT have any statistically significant relationship with the final translation quality or productivity when doing MTPE tasks?* To answer this question from a novel point of view, we conducted a human-computer interaction-informed study, where we recruited 11 professional translators in the English-Spanish legal translation com-

bination and carried out a pre-task questionnaire to examine their opinions, past experiences and attitudes towards MT. This questionnaire was followed by the translation of 13 texts using an interactive MT workflow. Subsequently, a professional, expert reviewer assessed the quality of the translations after ensuring consistent evaluation criteria with three professional reviewers. We examined the data obtained using different statistical analysis methods to find out whether there was any correlation between the past experiences and attitudes of translators towards MT and their resulting translation quality and productivity. Our hypothesis is that translators with negative pre-task perceptions of MT may produce translations with lower quality than their peers with positive pre-task perceptions because their predisposal to interact with MT will affect their translation processes. The following sub-sections describe the methodology used in-depth.

3.1 Participants

We recruited 11 professional translators in the English-Spanish language combination at an hourly rate of €20. To do this, we posted a job advert on ProZ (one of the most prominent job search platforms in the language services world) and X (which also has a large translator community). By posting on two different platforms and hiring participants on a first-come, first-served basis, we wanted to reach a large number of people without introducing any bias in the selection of participants. Participants were hired as long as they met the three basic conditions for participation: i) be native Spanish translators, ii) have professional experience in legal translation, and iii) have less than 5 years of professional experience. We decided to include the experience limitation because we wanted to minimise bias due to variable levels of experience. In addition, the translators were to perform post-editing, and previous studies suggested that people with more years of experience tended to have more problems interacting with technologies and were more likely to reject their daily use (Alabau et al., 2016).

Translators performed thirteen post-editing sessions of 45 minutes in Lilt (Green, 2016) over ten consecutive days (two weeks). In these sessions, three sessions were conducted through traditional post-editing, and ten sessions through interactive post-editing. The tasks were divided this

way for reasons of the project in which the present study is framed, but this has no impact on the data shown here, as all translators worked with the same texts, under the same conditions and had the same amount of time to translate. Translators were instructed to “Perform a full post-editing of the text, with the goal of achieving a perfectly fluent and adequate translation for a client in the legal domain. Any mistranslation may have critical legal consequences for the client, so ensure that you offer a professional translation. There is no problem if you do not finish the whole text in the allocated time”.

3.2 Translators’ pre-task perceptions

In order to collect translators’ pre-task perceptions of MT, we created an online questionnaire to be completed before starting the post-editing task. This included the following questions.

- Experience in MTPE tasks: How long have you engaged with MTPE tasks? Give an approximate time of use with months or years and months (e.g., 1 year and 6 months). [These experiences were then normalized to the number of months].
- Do you like MTPE?: On a scale of 1-7, where 1 is “Strongly Dislike” and 7 is “Strongly Like”, please rate your perception of doing MTPE tasks in professional translation projects.
- Do you trust MTPE?: On a scale of 1-7, where 1 is “Not trustworthy at all” and 7 is “Very trustworthy”, please rate if you can trust MTPE to help you successfully delivery a professional translation project.
- MT as a threat: Please rate how much you agree or disagree with this statement: “Machine Translation is a threat to the sustainability of the translation profession (Score 1 is “Disagree”, Score 7 is “Agree”).
- Is MTPE boring?: Please rate the following statement: “When I am doing MTPE tasks, I find them [SCORE]”. (Score 1 is “Boring”, Score 7 is “Engaging”).

The responses to the questionnaire were the translators’ pre-task perceptions that we correlated with final translation quality and productivity to examine if there was any relationship between them.

3.3 Translation Quality Evaluation

We worked with legal contracts in the English-Spanish combination and controlled the difficulty and length of the texts so that all translators worked with thirteen equally complex texts. For each text, a new task was set, and no translation memory was added. Difficulty was controlled with the Flesch-Kincaid index and the type-token ratio (Graesser et al., 2004). The total number of words translated and evaluated were 120,102. After obtaining the translations, translation quality was evaluated via human evaluation by using 1-4 Adequacy and Fluency scores with a professional, expert evaluator.

Although there are many different methods for evaluating translation quality (Moorkens et al., 2018a; Drugan, 2013), to answer our research question (*Do translators' (positive or negative) pre-task perceptions of MT have any statistically significant relationship with the final translation quality or productivity when doing MTPE tasks?*), we needed to obtain a final score of the translation quality of each translator. We considered the Adequacy and Fluency assessment to be the most appropriate method for our study, as it allowed us to obtain very detailed quality scores for each translator from two different points of views and has been extensively used in MT evaluation (Kocmi et al., 2022; Barrault et al., 2020). We discarded the MQM-based assessment (Freitag et al., 2021) because it focuses on the precise types of errors, and we did not need such a granular translation quality evaluation, plus it increased substantially the translation quality evaluation costs.

Best practices in human evaluation of translation quality recommend using several evaluators to reduce any potential subjectivity (Freitag et al., 2021). As an alternative, we have decided to follow common best practices in the fields of Computer Science and Information Retrieval (Artstein, 2017), also with recognised and widely-accepted methods for reducing evaluator subjectivity, and we have implemented the evaluation only with one expert reviewer after refining the evaluation criteria with a total of three reviewers through two different iterations. The scoring guidelines were updated after each iteration. The three reviewers were recruited by following the same methodology used for recruiting the translators, which can be found in the section 3.1 above, and they had +5 years of professional experience. The process followed for the quality evaluation was as follows:

First, we created a document explaining in detail the quality evaluation task to be carried out. Detailed scoring guidelines were also designed, in which each possible score (both for Adequacy and Fluency) was described in detail, and two examples were included for each type of score. The aim of these guidelines was to homogenise the evaluation criteria, and thus make the study and the results reproducible and reliable, trying to reduce the personal and subjective bias of each evaluator.

Once the first draft of the scoring guidelines was devised (containing two examples for every type of Adequacy and Fluency mistake), 50 translated segments were sent to the three reviewers. Texts evaluated in the iterations were fragments of English-Spanish legal contracts, similar in content and difficulty to the bulk of translations. The three reviewers annotated the translations and evaluated them according to the criteria of the scoring guidelines.

Subsequently, the Inter-Annotator Agreement (IAA) was calculated using Fleiss' Kappa. IAA can range between 0 and 1 and, generally, an IAA above 0.8 indicates that the consistency between annotations is high (Artstein, 2017). The IAA for our first round of annotations (Iteration 1) was 0.83, indicating that the scoring guidelines were clear, that the annotation consistency of the evaluators was high, but that there was still room for improvement.

Then, a Zoom meeting was held with the 3 reviewers to discuss the discrepancies of annotation in Iteration 1, and the scoring guidelines were updated with additional examples after some discussion. The main changes included re-wording and clarifying the annotation criteria, and more detailed explanations of the annotation limits, with the aim of improving the homogeneity of the annotation and increasing the consistency of the evaluations. Iteration 2 was then prepared, with 50 new segments, to be annotated by following the updated scoring guidelines in the same procedure as in Iteration 1. The IAA of Iteration 2 increased to 0.95, reflecting that the second version of the guidelines was clearer and more concise, and that a consistent evaluation could be obtained when evaluating translations if the guidelines were followed¹.

We then evaluated all the translations (120,102

¹Link to the final scoring guidelines: <https://zenodo.org/records/11091928>

words; 13 translations per each of the 11 translators) with a single reviewer, who we considered as the expert reviewer after the first two iterations, the homogenisation of criteria, and the updating of the scoring guidelines. To corroborate that the expert reviewer was still maintaining the annotation criteria halfway through the evaluation of the texts, the other reviewers, who participated in the earlier iterations, performed a cross-check evaluation. For this cross-check evaluation, 250 segments were randomly selected for annotation and the level of consistency was recalculated. The resulting IAA from the cross-check was 0.88, which also indicated a high consistency in the annotation criteria according to the elaborated scoring guidelines.

The end result is an Adequacy and Fluency score for each translator at the segment level. However, after translators performed the post-editing tasks, we observed that, in the allocated time for translation, some translators finished the texts, while others did not. Thus, the Adequacy and Fluency results have been normalised by calculating the average of all the segments translated by each translator. This normalisation has been carried out independently for both Adequacy and Fluency. By doing this, we can compare the results without any bias and independently for Adequacy and Fluency. Thus, we have a global quality score for each translator, ranging from 1 to 4².

3.4 Translation Productivity

Translation productivity was tracked in the CAT tool through a word per hour (WPH) measurement. In other words, we collected the translation productivity of every translator in each of the texts in WPH.

3.5 Statistical Analyses

First, we plotted every variable (translators' pre-task perceptions, fluency scores, adequacy scores and productivity measurements) in histograms to see whether the variables were normally distributed, and to strengthen our methodology we also performed the Shapiro-Wilk's test. As data violated the assumptions of normal distribution (p over .05), we conducted a Kendall's T correlation test for all the variables so as to explore

²The dataset including the source texts, the translations, the quality scores for fluency and adequacy, as well as the productivity measures can be found in: <https://zenodo.org/records/11092027>.

the relationships between the measures collected (Mellinger and Hanson, 2016). Due to the number of correlations performed increasing the likelihood of type I error, we recommend interpreting correlations at the .05 level with caution. In addition, it is worth stressing that the strength of the correlation coefficients vary according to the statistical test conducted. Therefore, by following Schober, Boer, and Schwarte (2018) advice, we interpret Kendall T's correlation coefficient strength in the following form: Weak (0.06-0.25), Moderate (0.26 to 0.49), Strong (0.50 to 0.71), and Very strong (0.71 to 1). Below, different heatmaps display the correlation coefficients of every pre-task perception variable in relationship to adequacy, fluency, and productivity. Variables containing an asterisk "*" indicate a statistically significant correlation. Also, the p-values are given for every variable in the wording.

4 Results

This section presents the correlations of translators' pre-task perceptions with fluency, adequacy and productivity.

4.1 Translators' pre-task perceptions of MT and fluency

Figure 1 shows in a heatmap the correlation coefficients resulting from the statistical analysis by considering translators' pre-task perceptions and final fluency scores for each of the texts. By looking at Figure 1, we can see that translators' feeling of boredom or engagement when performing MTPE assignments in a professional environment ($r(10) = -.012$, $p = .85$) showed no statistically significant correlation with Fluency scores.

However, all the other pre-task perceptions variables showed a statistically significant correlation with Fluency. On the one hand, we can observe two variables that show statistically significantly weak correlations. Whether translators had more or less experience in conducting MTPE tasks had no particular relationship with fluency results ($r(10) = .12$, $p = .04$). This means that, even if translators were new to interacting with MT in a professional environment, their fluency was not different to those translators with experience in post-editing. In a similar way, translators' attitude towards liking or disliking post-editing tasks ($r(10) = -.21$, $p = .0007$) showed a weak statistically significant correlation; this means that we cannot claim a re-

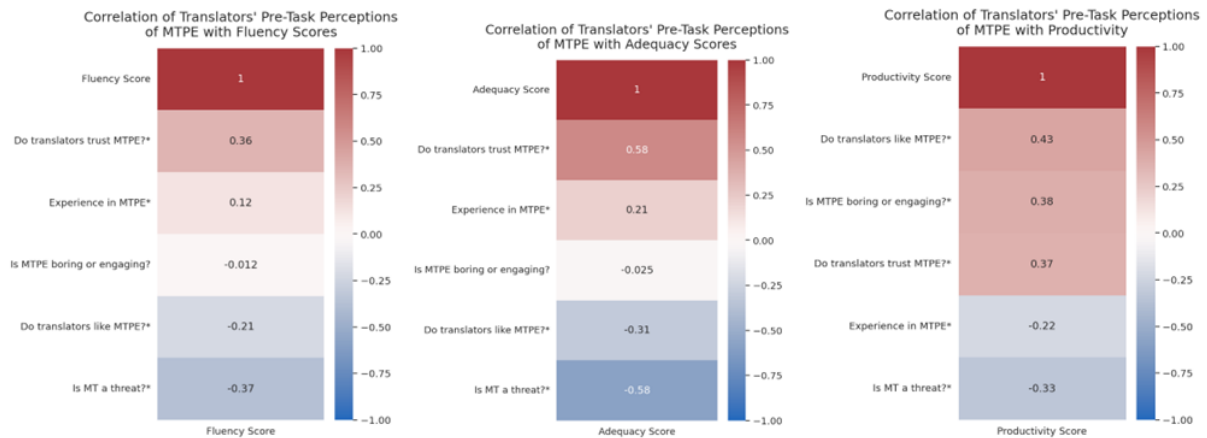


Figure 1: Correlation of translators' pre-task perceptions of MTPE with fluency, adequacy and productivity

relationship between the fact that translators liked post-editing or not with the production of more fluent translations. On the other hand, we can observe two variables that show stronger relationships. Translators' pre-task perceptions of MT being a threat to the translation profession ($r(10) = -.37, p = .0001$) and the level of trust they had on MTPE ($r(10) = .36, p = .0001$) showed statistically significant moderate correlations. In other words, this means that translators who had higher levels of trust in MTPE tasks as an aid in their professional translation projects tended to report higher fluency scores. In a similar way, those translators who thought that MT was a threat to their profession tended to produce less fluent translations.

4.2 Translators' pre-task perceptions of MT and adequacy

The correlation heatmap of Figure 1 provides a visual summary of the statistical analysis conducted on translators' pre-task perceptions of MTPE and their translation quality results, specifically focusing on translation adequacy.

One of the most notable results is that we observed a strong statistically significant positive correlation ($r(10) = .58, p = 0.0001$) between translators' trust in MTPE and the adequacy scores, implying that higher trust in the system is linked to higher performance levels in producing adequate translations. The other remarkable result is that translators' view of MT as a threat yielded a strong statistically significant negative correlation ($r(10) = -.58, p = 0.0001$), suggesting that apprehensions about the technology's impact on the profession may undermine translation adequacy.

Factors such as the enjoyment of conducting

MTPE tasks ($r(10) = -.31, p = 0.0001$), and the overall experience in MTPE ($r(10) = .21, p = 0.0004$), showed less pronounced yet statistically significant correlations, indicating that these perceptions might not be as critical in influencing the adequacy of translation outcomes.

These insights contribute to the ongoing discourse on the human factors influencing contemporary translator-computer interactions, underscoring the complex interplay between subjective perceptions and objective translation performance metrics. These results indicate that translators' lack of trust in MTPE tasks and the consideration of MT as a threat to their profession may have a strong effect on translation quality, especially adequacy, even before the task has already started.

4.3 Translators' pre-task perceptions of MT and productivity

In terms of productivity, Figure 1 provides a quantitative depiction of the correlations between translators' pre-task perceptions of MTPE and their measured productivity in WPH. Here, the results from every pre-task perception variable were statistically significant.

The data indicates that positive perceptions towards MT, such as liking MTPE tasks as a professional aid ($r(10) = .43, p = 0.0001$), finding them engaging ($r(10) = .38, p = 0.0001$), or the level of trust in MT ($r(10) = .37, p = 0.0001$) are moderately correlated with higher productivity scores. Conversely, the negative moderate correlation with the perception of MT as a professional threat ($r(10) = -.33, p = 0.0001$), although displaying a weaker association, highlights potential areas of concern. These findings may reflect a complex-

ity in MTPE's perceived impact on the translation industry, which could influence translator productivity.

There is also a weak statistically significant negative correlation with MTPE experience ($r(10) = -.22, p = 0.0001$).

5 Discussion of the results

After analysing the results, we can see that translators' pre-task perceptions of MT have a higher correlation with adequacy scores than with fluency and productivity scores. Research by Castilho et al. (2017) suggested that NMT systems produce very fluent translations and, therefore, the effect of translators' pre-task perceptions may not be that impactful on fluency scores if the system already offers a fluent MT output. In terms of adequacy scores, however, translators' pre-task perceptions have a bigger impact. These results indicate that, in post-editing tasks with NMT systems that offer good MT quality (English-Spanish in our case), adequacy scores have a higher dependence on the translator, while fluency scores have a lower dependence on translators' translation and/or post-editing skills because the MT system already offers higher quality MT output. The expert reviewer assessed the MT output of the MT system used in this study, which obtained a global score of 3.48/4 in terms of adequacy and 3.71/4 in terms of fluency, showcasing good quality results. It may be interesting to further validate this idea by replicating this study in a different language combination, particularly in a case in which NMT systems offer lower MT quality (i.e. a low-resource language combination). In terms of productivity, in general terms, we can see that translators who had positive pre-task perceptions of MT tended to translate faster when conducting post-editing tasks than their peers with negative pre-task perceptions of MT.

Experience in MTPE tasks has no correlation with final translation quality or productivity in the data analysed. These results indicate that it is the translator who matters. Results suggest that a professional translator with good translation skills will offer good quality translations and will work equally faster when interacting with MT, regardless of their experience in providing MTPE language solutions.

It is worth stressing that the most notable correlation coefficients were observed in two specific

variables: translators' level of trust in MTPE and the perception of MT being a threat to the translation profession. The level of trust translators have in MT shows a strong correlation with adequacy and a moderate correlation with fluency and productivity. This is interesting, as translators who trust MT systems to help them work in their professional, daily tasks offered higher final translation quality than those who did not trust MT systems. This is in line with previous research in cognitive science (Albarracín, 2021), which indicated that prior negative perceptions are an important and crucial determinant for future attitudes and behaviours. In our case, translators' pre-task perceptions of MTPE tasks had a strong negative correlation with the quality of the final product, that is, the translation. This may be because translators who do not trust MT do not enjoy this interaction or do not give their best when interacting with MT in their regular, professional workflows. This also applies to productivity: translators with higher levels of trust on MT translated faster, probably because they were more enthusiastic about engaging with MT. Those translators who did not trust MT were probably more reluctant to engage with MT in the best of their abilities.

This backs up the results of the second pre-task perception variable with a strong negative correlation in our study, that is, whether translators consider MT as a threat to the translation profession. The perception of MT being a threat to translators showed a strong association with final translation quality, both in terms of Adequacy and Fluency, and the results showed statistical significance. The correlation of this pre-task perception variable with productivity is weaker, but still moderate and statistically significant. What these correlations mean is that translators who, even before starting a post-editing task, think that MT is a threat and harmful for the translation profession are more likely to produce lower translation quality and to translate slower.

The study's findings on the relationship between translators' pre-task perceptions and translation quality and productivity have profound implications, offering novel insights into the dynamics of modern translator-computer interactions. It is evident that the approach translators adopt towards a task plays a critical role in determining the final outcome, with varying degrees of influence on different aspects of translation performance (in terms

of quality or productivity). These results highlight the problem of modern translator-computer interactions, and suggest that we should give higher attention to the improvement of these interactions, probably by looking at the MTUX (Koponen et al., 2020; Karakanta et al., 2022; Briva-Iglesias and O'Brien, 2023), putting the human in the centre of modern translator-computer interactions (Shneiderman, 2022), reducing recent complaints about dehumanisation (Moorkens, 2020). Also, this has great implications for the training of translators, as it highlights the importance of MT literacy from different points of view (Bowker and Ciro, 2019), which are further detailed in the conclusions.

6 Conclusions

Since the emergence of language technologies, the perceptions of those who interact with them have been studied, ranging from professional translators to gisting users (Nurminen, 2019). Particular attention has been paid to the perceptions of professional translators, who saw their traditional workflows disrupted by these new technologies. However, most studies to date have been descriptive and did not take into account the relationship between perceptions and the quality of the translation or the productivity of the translator. This article aims to fill this gap in the literature with a longitudinal study of 11 professional English-Spanish translators in the legal domain, to explore whether negative or positive translators' pre-task perceptions of MT have any relationship with the final translation quality and productivity in a post-editing workflow. In terms of limitations, it would have been better to increase the sample size of the translations or the number of translators. However, we hired 11 professional translators who produced a total of 120,102 words over 10 consecutive days. It would have also been ideal to have the three reviewers assess all the translations, but due to budget constraints, we reduced the reviewer bias through different evaluation iterations, the measurement of IAA, and the refinement of a set of quality scoring guidelines. Exploring additional domains to the legal field would have also been positive to assess whether these results are generalisable to other translation specialisations.

As a conclusion, the results suggest that translators with negative pre-task perceptions of MT tend to deliver poorer quality translations, as well as to translate slower, than their peers with positive

pre-task perceptions of MT. Specially, in our study, we observed that translators who thought that MT was a threat to their profession or distrusted MT as an aid in their work obtained lower quality and productivity scores. These were the two variables with the highest correlation coefficients. By contrast, translators with positive pre-task perceptions obtained better translation quality and productivity scores. This may have happened because translators with negative pre-task perceptions of MT may have not interacted with MT adequately or with an open-minded point of view, impacting their final translation quality and productivity.

This research opens up new questions: do these results suggest that there is a direct relationship between pre-task perceptions of translation technologies and final quality and productivity results? Would translators with negative pre-task perceptions of MT obtain better quality results if they translated without MT? And what would happen if we trained these translators and taught them to see MT as an aid to augment their skills and help them in their professional tasks?

Although we now have new questions to answer, what is clear is that the level of trust in MT and the conception of MT as a threat to the translation profession have a strong correlation with final quality results, especially Adequacy, and a moderate correlation with productivity. These correlations have important implications. Translators fearing MT, or those who are more reluctant to trust their interactions with MT, may not be leveraging the advantages and benefits MT offers them. Therefore, the results call for multiple actions to be taken in order to:

- Increase translators' confidence in their interactions with MT as a tool that can be useful and support them in professional projects, always bearing in mind that translators are the ones controlling the interaction, and that MT functions as a support that can offer alternative terminology solutions or facilitate understanding of the source text, among other forms of assistance. The main goal of technologies should be to augment translators and reduce their human cognitive limitations (O'Brien, 2023; Raisamo et al., 2019; Alicea, 2018; Shneiderman, 2022), pursuing human-centered, augmented machine translation (Briva-Iglesias, 2024), not to replace and substitute them.

- Present MT as a tool that can facilitate translators' work, either to increase productivity or to open doors to new professional markets and domains, as is the case of *language engineers* (Briva-Iglesias and O'Brien, 2022). As Stasimioti and Sosoni (2019) reported, training translators on MT will change their perceptions of MT because they will learn what MT allows them to do or not. Translators' technological and MT literacy is now more important than ever in the AI age (Bowker and Ciro, 2019).

These two elements would increase the adoption and use of MT as assistance, as well as reduce the negativity of translators' pre-task perceptions of MT. However, it is vital to stress that this MT literacy must be accompanied by a broad and holistic view of MT, including its limitations, so that translators acquire a critical view of when it is appropriate and when not to use MT, as it may also involve important ethical issues (Moorkens, 2022).

Acknowledgements

This work was conducted with the financial support of the Science Foundation Ireland Centre for Research Training in Digitally-Enhanced Reality (d real) under Grant No. 18/CRT/6224.

References

- Alabau, Vicent, Michael Carl, Mercedes García-Martínez, and Jesús González-Rubio. 2016. *Learning Advanced Post-Editing*.
- Albarracín, Dolores and Robert S. Wyer. 2000. The Cognitive Impact of Past Behavior: Influences on Beliefs, Attitudes, and Future Behavioral Decisions. *Journal of personality and social psychology*, 79(1):5–22, July.
- Albarracín, Dolores, editor. 2021. *The Impact of Past Experience and Past Behavior on Attitudes and Behavior*. Cambridge University Press.
- Alicea, Bradly. 2018. *An Integrative Introduction to Human Augmentation Science*.
- Artstein, Ron. 2017. *Inter-Annotator Agreement*. Springer Netherlands.
- Barrault, Loïc, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. Findings of the 2020 Conference on Machine Translation (WMT20). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1–55. Association for Computational Linguistics.
- Bender, Emily M., Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? . In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623. Association for Computing Machinery.
- Bowker, Lynne and Jairo Buitrago Ciro. 2019. *Machine Translation and Global Research: Towards Improved Machine Translation Literacy in the Scholarly Community*. Emerald Publishing Limited.
- Briva-Iglesias, Vicent and Sharon O'Brien. 2022. The Language Engineer: A Transversal, Emerging Role for the Automation Age. *Quaderns de Filologia - Estudis Lingüístics*, 27:17–48.
- Briva-Iglesias, Vicent and Sharon O'Brien. 2023. Measuring Machine Translation User Experience: A Comparison between AttrakDiff and User Experience Questionnaire. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 335–344.
- Briva-Iglesias, Vicent, Sharon O'Brien, and Benjamin R. Cowan. 2023. The impact of traditional and interactive post-editing on Machine Translation User Experience, quality, and productivity: *Translation, Cognition & Behavior*, 6(1).
- Briva-Iglesias, Vicent. 2024. *Fostering human-centered, augmented machine translation: analysing interactive post-editing*. PhD thesis. Dublin City University.
- Brown, Tom B., Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners.
- Cadwell, Patrick, Sharon O'Brien, and Carlos S. C. Teixeira. 2018. Resistance and accommodation: Factors for the (non-) adoption of machine translation among professional translators. *Perspectives*, 26(3):301–321.
- Castilho, Sheila, Joss Moorkens, Federico Gaspari, Iacer Calixto, John Tinsley, and Andy Way. 2017.

- Is Neural Machine Translation the New State of the Art? *The Prague Bulletin of Mathematical Linguistics*, 108(1):109–120.
- de Almeida, Giselle. 2013. *Translating the Post-Editor: An Investigation of Post-Editing Changes and Correlations with Professional Experience across Two Romance Languages*. Doctoral, Dublin City University.
- Drugan, Joanna. 2013. *Quality in Professional Translation: Assessment and Improvement*, volume 9. A&C Black.
- ELIS Research. 2023. EUROPEAN LANGUAGE INDUSTRY SURVEY 2023.
- Etchegoyhen, Thierry, Anna Fernández Torné, Andoni Azpeitia, Eva Martínez Garcia, and Anna Matamala. 2018. Evaluating Domain Adaptation for Machine Translation Across Scenarios. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*. European Language Resources Association (ELRA).
- Firat, Gökhan. 2021. Uberization of translation: Impacts on working conditions. *The Journal of Internationalization and Localization*, 8(1):48–75.
- Freitag, Markus, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, Errors, and Context: A Large-Scale Study of Human Evaluation for Machine Translation. *arXiv:2104.14478 [cs]*, April.
- Gaspari, Federico, Antonio Toral, Sudip Kumar Naskar, Declan Groves, and Andy Way. 2014. Perception vs reality: Measuring machine translation post-editing productivity. In *Third Workshop on Post-Editing Technology and Practice*, volume 60.
- Graesser, Arthur C., Danielle S. McNamara, Max M. Louwerse, and Zhiqiang Cai. 2004. Coh-Metrix: Analysis of text on cohesion and language. *Behavior Research Methods, Instruments, & Computers*, 36(2):193–202.
- Green, Spence. 2016. Interactive Machine Translation. In *Conferences of the Association for Machine Translation in the Americas*, page 93.
- Guerberof Arenas, Ana, Joss Moorkens, and Sharon O’Brien. 2021. The impact of translation modality on user experience: An eye-tracking study of the Microsoft Word user interface. *Machine Translation*, 35(2):205–237.
- Karakanta, Alina, Luisa Bentivogli, Mauro Cettolo, Matteo Negri, and Marco Turchi. 2022. Post-editing in Automatic Subtitling: A Subtitlers’ perspective. In *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, pages 261–270, Ghent, Belgium. European Association for Machine Translation.
- Kocmi, Tom, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 Conference on Machine Translation (WMT22). In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costajussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45. Association for Computational Linguistics.
- Koponen, Maarit, Umut Sulubacak, Kaisa Vitikainen, and Jörg Tiedemann. 2020. MT for Subtitling: Investigating professional translators’ user experience and feedback. In Ortega, John E., Marcello Federico, Constantin Orasan, and Maja Popovic, editors, *Proceedings of 1st Workshop on Post-Editing in Modern-Day Translation*, pages 79–92. Association for Machine Translation in the Americas.
- Matusov, Evgeny, Patrick Wilken, and Yota Georgakopoulou. 2019. Customizing Neural Machine Translation for Subtitling. In *Proceedings of the Fourth Conference on Machine Translation (Volume 1: Research Papers)*, pages 82–93. Association for Computational Linguistics.
- Mellinger, Christopher and Thomas Hanson. 2016. *Quantitative Research Methods in Translation and Interpreting Studies*. Routledge.
- Moorkens, Joss, Sheila Castilho, Federico Gaspari, and Stephen Doherty, editors. 2018a. *Translation Quality Assessment: From Principles to Practice*, volume 1 of *Machine Translation: Technologies and Applications*. Springer International Publishing.
- Moorkens, Joss, Antonio Toral, Sheila Castilho, and Andy Way. 2018b. Translators’ perceptions of literary post-editing using statistical and neural machine translation. *Translation Spaces*, 7(2):240–262.
- Moorkens, Joss. 2020. “A tiny cog in a large machine”: Digital Taylorism in the translation industry. *Translation Spaces*, 9(1):12–34.
- Moorkens, Joss. 2022. Ethics and machine translation. *Machine translation for everyone*, pages 121–140.
- Nurminen, Mary. 2019. Decision-making, Risk, and Gist Machine Translation in the Work of Patent Professionals. In *Proceedings of the 8th Workshop on Patent and Scientific Literature Translation*, pages 32–42. European Association for Machine Translation.

- O'Brien, Sharon, Maureen Ehrensberger-Dow, Marcel Hasler, and Megan Connolly. 2017. Irritating CAT tool features that matter to translators. *Hermes: Journal of Language and Communication in Business*, 56:145–162.
- O'Brien, Sharon. 2023. Human-Centered augmented translation: Against antagonistic dualisms. *Perspectives*, pages 1–16.
- Olohan, Maeve. 2011. Translators and translation technology: The dance of agency. *Translation Studies*, 4(3):342–357.
- Pérez-Macías, Lorena, María del Mar Sánchez Ramos, and Celia Rico. 2020. Study on the Usefulness of Machine Translation in the Migratory Context: Analysis of Translators' Perceptions. *Open Linguistics*, 6(1):68–76.
- Piller, Ingrid, Jie Zhang, and Jia Li. 2020. Linguistic diversity in a time of crisis: Language challenges of the COVID-19 pandemic. *Multilingua*, 39(5):503–515.
- Raisamo, Roope, Ismo Rakkolainen, Päivi Majaranta, Katri Salminen, Jussi Rantala, and Ahmed Farooq. 2019. Human augmentation: Past, present and future. *International Journal of Human-Computer Studies*, 131:131–143.
- Rossi, Caroline and Alice Carré. 2022. How to choose a suitable NMT solution?: Evaluation of MT quality.
- Rossi, Caroline and Jean-Pierre Chevrot. 2019. Uses and perceptions of Machine Translation at the European Commission. *The Journal of specialised translation (JoSTrans)*.
- Schober, Patrick, Christa Boer, and Lothar A. Schwarte. 2018. Correlation Coefficients: Appropriate Use and Interpretation. *Anesthesia & Analgesia*, 126(5):1763.
- Shneiderman, Ben. 2022. *Human-Centered AI*. Oxford University Press.
- Stasimioti, Maria and Vilemini Sisoni. 2019. Undergraduate Translation Students' Performance and Attitude vis-à-vis Machine Translation and Post-editing: Does Training Play a Role. In *41st Translating and the Computer Conference*, pages 125–136.

Bayesian Hierarchical Modelling for Analysing the Effect of Speech Synthesis on Post-Editing Machine Translation

Miguel Rios, Justus Brockmann, Claudia Wiesinger,
Raluca Chereji, Alina Secară, Dragoş Ciobanu

Centre for Translation Studies, University of Vienna

{miguel.angel.rios.gaona, justus.brockmann,
claudia.wiesinger, raluca-maria.chereji,
alina.secara, dragos.ioan.ciobanu}@univie.ac.at

Abstract

Automatic speech synthesis has seen rapid development and integration in domains as diverse as accessibility services, translation, or language learning platforms. We analyse its integration in a post-editing machine translation (PEMT) environment and the effect this has on quality, productivity, and cognitive effort. We use Bayesian hierarchical modelling to analyse eye-tracking, time-tracking, and error annotation data resulting from an experiment involving 21 professional translators post-editing from English into German in a customised cloud-based CAT environment and listening to the source and/or target texts via speech synthesis. We find that using speech synthesis in the PEMT task has a non-substantial positive effect on quality, a substantial negative effect on productivity, and a substantial negative effect on the cognitive effort expended on the target text, signifying that participants need to allocate less cognitive effort to the target text.

1 Introduction

The growing adoption of data-driven approaches to machine translation (MT) since the 2000s (Kenny, 2020) has brought ongoing change to the practice of translation. While ‘standard’ human translation still appears to be the dominant type of service, industry surveys have repeatedly identified post-editing of MT (PEMT) as the service with the highest growth potential, according to language service

providers (ELIA et al., 2023). A wealth of previous research has addressed the implications of this change, ranging from potential productivity gains (Plitt and Masselot, 2010; Lüubli et al., 2019) to impacts on creativity (Guerberof-Arenas and Toral, 2022). A central theme in studies on PEMT is the effort expended by translators (Krings, 2001) and how it might be impacted by the tools they use. Moreover, previous work has probed how well PEMT is supported by the user interfaces used by translators (Moorkens and O’Brien, 2017; Herbig et al., 2020), indicating room for improvement.

A relatively novel approach to supporting PEMT processes – and translation in general – is integrating automatic text-to-speech synthesis (Taylor, 2009) in computer-assisted translation (CAT) tools. The idea is for the translator to be able to trigger an artificial voice that ‘reads’ to them the source and/or target text, thus adding a new mode of text reception to information processing approaches that have traditionally relied heavily on reading. Only little attention has thus far been given to this method in related work, but initial findings point to potential benefits in revision (Ciobanu et al., 2019) and PEMT (Wiesinger et al., 2022). This motivates our present study into the impact of speech synthesis on the PEMT process.

In this paper, we measure the effect of adding text-to-speech into a translation workflow for PEMT for the English-German language pair. We focus on the the target text quality delivered, cognitive effort expended, and productivity recorded, with an emphasis on the statistical modelling approach. Eye-tracking output metrics, such as the number or duration of fixations on both source and target segments are used to measure the cognitive effort during PEMT (Moorkens, 2018). Moreover, linear models and linear mixed effect models are

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

commonly used for the analysis of eye-tracking studies (Kim et al., 2022; Silva et al., 2022), including the use of linear models to investigate the relation across text complexity of the source, cognitive effort, and PEMT (Dai and Liu, 2024). Silva et al. (2022) discuss the disadvantages of standard statistical tests for eye-tracking data analysis in subtitling. For example, t-tests conflate the data by averaging a participant into one outcome variable, and ignore other variables (predictors) that may affect the results of an experiment. Instead, they use a linear mixed-effect model for the analysis of cognitive effort (outcome variable) of reading subtitles given the effect of subtitle speeds. However, linear models require large amounts of data to achieve reliable learned estimates (Silva et al., 2022). Bayesian hierarchical models cope with data scarcity by adding information from the data structure, and prior expert knowledge that works as a regulariser to avoid over-fitting to the available data (Gelman et al., 2004). We use Bayesian hierarchical modelling to tackle the issue of data scarcity that is common in eye-tracking studies (O’Brien, 2009). Our contributions are as follows:

- We report on a PEMT study including eye-tracking, time-tracking, and quality evaluation.
- We introduce a Bayesian hierarchical model for PEMT data analysis.
- We measure how a speech-enabled mode of working may support professional translators post-editing within a CAT tool.

2 Methodology

2.1 Participants

The participants were recruited via the network of the language service provider Translated, the professional translator association UNIVERSITAS Austria, the Austrian Economic Chambers (WKO), and the website of the HAITrans research group¹. Prospective participants were asked to fill in a recruitment questionnaire to determine whether they fulfilled the participation requirements. In total, we recruited 21 professional translators working from English into German who have German as a first language. All translators have at least three years of professional translation experience, with 10 participants having over 11 years of experience.

¹<https://haitrans.univie.ac.at/>

Most participants have at least one year of PEMT experience, although five translators have little to no PEMT experience. The experiment received ethical approval from the Ethics Committee of the University of Vienna. All participants were remunerated for their time. After the conclusion of the experiment, the participant data were anonymised, and the participants were assigned an experiment ID.

2.2 Materials

The source texts used in the experiment consisted of four excerpts from two separate factsheets produced by the International Federation of Red Cross and Red Crescent Societies, UNICEF, and the World Health Organisation about stigma, mistrust, and denial in relation to COVID-19. Both factsheets were published online on the British Red Cross’s Community Engagement Hub² in 2020.

The four English source text parts have a combined total number of 1,423 words, with their respective IDs being text 1 (t1), text 2 (t2), text 3 (t3), and text 4 (t4). To counteract the impact of the text parts on the results, we alternated text 2 and text 3 for every other participant. For this reason, we ensured comparability of the four text parts in terms of standard measurements of linguistic complexity and lexical richness as shown in Table 1, as well as readability as shown in Table 2. We use Textstat³ for the readability scores, and LexicalRichness⁴ for the linguistic complexity and lexical richness scores. The Flesch–Kincaid Reading Ease scores class all text IDs as fairly easy to read (between 80.0-70.0) and at 7th grade level. All text IDs have a consistent low linguistic complexity expressed as Type-Token Ratio (TTR).

Text	Word count	Number of syllables	Standardised TTR	Sentence count	Average sentence length
t1	342	454	0.483	18	19.0
t2	374	498	0.475	18	20.8
t3	352	471	0.520	18	19.6
t4	355	477	0.532	19	18.7

Table 1: Linguistic complexity and lexical richness for each text ID.

²<https://communityengagementhub.org/>

³<https://github.com/textstat/textstat>

⁴<https://github.com/lsys/lexicalrichness>

Text	Flesch Reading Ease	Flesch-Kincaid Grade Level	New Dale-Chall
t1	77.57	7.2	7.58
t2	75.74	7.9	7.92
t3	76.96	7.4	7.75
t4	77.87	7.0	7.90

Table 2: Readability scores for each text ID.

2.3 Design

Before coming to the eye-tracking lab, the participants received a translation brief in German⁵ with information about the task scope, target audience and style requirements, as well as the requirements for PEMT. Those five participants with little to no prior post-editing experience were also sent a short training video on MT and PEMT to watch ahead of the experiment. Upon arrival, participants signed a declaration of consent, then filled in a pre-experiment questionnaire designed to collect some demographic information and to determine their exposure to CAT tools.

The participants' task in this experiment was to post-edit the four source text parts from English into German in a customised version of the CAT tool Matecat⁶ enhanced by Translated⁷ with a proprietary speech synthesis function. Participants worked in two conditions: in silence, and in a sound condition whereby they could trigger speech synthesis for the source and target segments.

An EyeLink Portable Duo eye tracker⁸ was used to record the participants' gaze during the experiment. Prior to performing these tasks, participants post-edited a short practice text using speech synthesis to familiarise themselves with the task setup and working environment. Each participant's computer screen and computer interactions were recorded for later annotation and comparison with other experiment participants. The total duration of the experiment was up to 3 hours.

2.4 Data Collection

The screen recordings, overlaid with participants' in-task gaze data captured with the eye tracker, were manually annotated in the SR Research Data Viewer software⁹. This included adding timestamps

⁵<https://github.com/HAITrans-lab/HAITrans-bayesian-multilevel-model>

⁶<https://www.matecat.com/>

⁷<https://imminent.translated.com/>

⁸<https://www.sr-research.com/eyelink-portable-duo/>

⁹<https://www.sr-research.com/data-viewer/>

for task start and end times and recording the number and type of exits from the Matecat environment (e.g., to look up terms online or read the source texts made available in Microsoft Word). Areas of interest were defined around the source and target text areas in Matecat to allow for using in the analysis only the gaze data that fell within these areas.

Reports containing measures such as the total number of fixations, dwell time, and mean fixation duration for the source and target sections of the video recordings, as well as the start and end timestamps of each trial, were then generated and used for the analysis. The post-edited target texts produced by the participants were exported from Matecat for subsequent annotation and quality evaluation by multiple contributors.

When conducting eye-tracking experiments, high participant attrition rates are to be expected (O'Brien, 2009). We were able to obtain eye-tracking measures for 19 out of the 21 participants. Furthermore, due to data corruption, data from *t1* is missing entirely for one of the participants. This explains the differences in participant numbers that can be seen in Tables 3, 5, 7, and 9.

2.5 Analysis

We use Bayesian hierarchical modelling for our data analysis (Gelman et al., 2004). Hierarchical models are also known as linear mixed effects models. The motivation to use Bayesian data analysis is the data scarcity (few observations), improved learned estimates, and uncertainty quantification of the estimates. Linear regression models learn the relation of a given measurement or outcome with one or multiple predictor variables (Gelman and Hill, 2007). For example, the positive or negative effect (linear relation) of the sound condition variable on the measured quality of the produced translations.

A hierarchical model outlines a hierarchy over the data where variables are considered related or grouped under the structure of a given problem (Gelman et al., 2004). Moreover, hierarchical models take advantage of their structure to improve the learned estimates by reducing variance when the data are limited. For example, we can define groups with the produced translations by participant, condition, or type of text. A hierarchical model consists of population-level effects (fixed) for variables that describe all the observed data, and group-level effects (random) for clusters or variables that describe

variability across groups (McElreath, 2016).

Bayesian linear models allow us to test the probability of our hypothesis given the observed data by providing a posterior distribution, which contains probable values of an effect. For uncertainty quantification, Bayesian linear models produce the credible interval (CI) that is a range containing a percentage of probable values (e.g. 95%). With the given data, the effect has 95% probability of falling within this range. Moreover, Bayesian models provide a posterior distribution for the learned estimates, instead of a point from standard regression models. The posterior distribution is used to analyse the direction and size of the effect, as well as the uncertainty.

The practical importance of an effect can be decided based on the region of practical equivalence (ROPE) (Kruschke, 2018). The ROPE is a range with a small or practically no effect, which is an area that encloses values that are equivalent to the null. As a decision rule, if a large part of an estimate 95% CI falls outside from the ROPE, the effect is considered **substantial** or of **practical importance** (Kruschke, 2018). The ROPE for linear models can be defined with the standard deviation (sd) of an outcome variable as $[-0.1 * sd(\text{outcome variable}), 0.1 * sd(\text{outcome variable})]$.

We are interested in analysing the following outcome variables Y : Quality score based on human error annotation, Productivity with words per hour (PEMT speed), Cognitive effort with the mean fixation duration on the source text (MFD-ST), and the mean fixation duration on the target text (MFD-TT).

For the predictor variables X , we use: Condition (no sound, and sound), ID of the text (t1, t2, t3, and t4), Number of external searches, and PEMT experience (yes, no). *Condition* refers to whether the participant used speech synthesis while post-editing (sound) or not (no sound). The *text ID* identifies the text part that was post-edited. The *number of external searches* specifies how many times the participant left the CAT tool interface to perform a web search or consult other sources. *PEMT experience* refers to a participant having (yes, y) or not (no, n) previous PEMT experience.

We define a hierarchical model with random intercepts and slopes. We use the participants as the second level grouping variable to measure the effect of the sound condition on each person, and the variability across them. The population-level effects are the X predictors, and intercept and slopes

for each condition and participant for group-level effects. The description of the hierarchical model is as follows:

$$\begin{aligned}
 y_i &\sim \mathcal{N}(\mu, \sigma^2) \\
 \mu &= \alpha_j[i] + \beta_{1j}[i](\text{condition}) \\
 &\quad + \beta_2(\text{text}) + \beta_3(\text{n_searches}) \\
 &\quad + \beta_4(\text{PEMT_experience}) \\
 \begin{pmatrix} \alpha_j \\ \beta_{1j} \end{pmatrix} &\sim \mathcal{N} \left(\begin{pmatrix} \mu_{\alpha_j} \\ \mu_{\beta_{1j}} \end{pmatrix}, \begin{pmatrix} \sigma_{\alpha_j}^2 & \rho_{\alpha_j \beta_{1j}} \\ \rho_{\beta_{1j} \alpha_j} & \sigma_{\beta_{1j}}^2 \end{pmatrix} \right) \\
 &\quad , \text{ for participant } j = 1, \dots, J
 \end{aligned}$$

where y_i is the outcome variable (e.g. quality score, PEMT speed) predicted from a normal distribution (regression) with mean μ based on a hierarchical linear model and variance σ^2 . For the linear model: α_i intercept and $\beta_{1...4}$ slopes with a uniform prior are population-level coefficients, α_j intercept and β_{1j} slopes are group-level coefficients with a normal prior for each participant j .

We use the brms package in R for our Bayesian analyses (Bürkner, 2017). brms provides an interface for Bayesian linear models, and hierarchical models using Stan¹⁰. We show the brms formulas for our hierarchical model in the Appendix A and the scripts for our experiments are available at: <https://github.com/HAITrans-lab/HAITrans-bayesian-multilevel-model>.

3 Results

3.1 Quality

To assess quality, we scored the post-edited texts using an error typology based on the Multidimensional Quality Metrics (MQM) framework (Burchardt, 2013). Two professional translators with more than three years of experience annotated the raw MT output for the four texts using the MQM typology within the CATMA annotation tool (Gius et al., 2023). These gold standard texts are labelled with all MT errors that the participants are expected to correct according to the translation brief. The annotators first labelled the texts independently of each other, and then combined their labels into the final gold standard, asking a third annotator for advice whenever they disagreed. The MQM error severities are defined with the following weights: Minor (1), Major (5), and Critical (25). To produce the quality score for each text, we counted the number of MT errors left uncorrected, as well as errors

¹⁰<https://cran.rstudio.com/web/packages/brms/>

newly introduced by our participants, and weighted them according to their severity. This resulted in a score between 0 and 100 for each text, where a score of 100 would mean there were no errors in the post-edited target texts.

Condition	Text	Variable	n	mean	sd
nos	t1	quality score	20	94.635	2.469
nos	t2	quality score	10	81.311	8.054
nos	t3	quality score	11	93.828	4.672
s	t2	quality score	11	86.922	7.839
s	t3	quality score	10	88.75	4.896
s	t4	quality score	21	94.271	2.528

Table 3: Summary statistics of the *quality score* with mean and standard deviation (sd).

Population-Level Effects			
Predictors	Estimate	CI (95%)	ROPE ↓
Intercept	96.85	[92.66, 101.20]	0.00%
condition [s]	0.36	[-2.12, 2.82]	41.00%
text [t2]	-10.64	[-13.41, -7.89]	0.00%
text [t3]	-3.51	[-6.22, -0.79]	0.00%
text [t4]	-0.66	[-4.09, 2.76]	30.00%
n searches	-0.13	[-0.51, 0.24]	100%
PEMT experience [y]	-2.16	[-6.55, 2.29]	15.00%
Group-Level Effects			
	sd	CI (95%)	
Intercept	3.85	[2.29, 5.85]	
condition [s]	1.00	[0.04, 2.82]	

Table 4: Summary of the fitted model for the *quality score*. ROPE size ± 0.66 .

Table 3 shows summary statistics with the number of participants (n), the mean, and sd of the quality score. We show the statistics grouped by both condition *no sound* (nos) and *sound* (s), and the ID of the text (t1, t2, t3, t4).

Table 4 shows the model summary for the quality score. The predictors for the population-level effects are summarised with estimate (learned mean), 95% credible interval (CI), and percentage of the estimate that overlaps with the ROPE. The linear model takes a class or name of a variable in alphabetical order as the reference for the Intercept and adds the value of the names left as the slopes. For example, the intercept is the no sound condition nos and the sound condition s is represented with the slope condition (s).

The sound condition has a non-substantial positive effect on the quality score, because the estimate 95% CI has a large overlap with the ROPE (41%).

To visualise the overlap of the sound condition CI with the ROPE, we refer the reader to Figure 6 in the Appendix. The texts t2, t3 have the highest substantial negative effect on the quality score. The effect of the number of searches (n searches) and having PEMT experience (y) are non-substantial. The group-level effect indicates how the condition (s) estimate varies from participant (group) to participant based on the sd.

To visualise the learned estimates, we show the conditional effects in Figure 1. The conditional effect plot shows the effects of each categorical or continuous predictor with the CI bar around the estimate on the outcome variable. In Figure 1 a) there is a large overlap between the CIs of the no-sound and sound conditions that indicates high uncertainty, and no difference between them. For Figure 1 b) the overlap for t2 between texts is little and indicates low uncertainty. Next, in Figure 1 c), a high number of external searches decreases the quality, but the uncertainty of the estimate is high. Moreover, in Figure 1 d), having PEMT experience (y) decreases the quality, but the difference compared to not having experience (n) is uncertain.

Figure 5 a) (Appendix) shows the fitted curve with the data points across texts from the quality score model. The posterior predictions plot shows the posterior mean (fit curve) and 95% credible interval (uncertainty bars) for each data point from the model. In other words, it plots the relation between each condition and the quality score. We can observe a difference in quality for t2, and under the sound condition, but it is small given the CI overlap.

3.2 Productivity

PEMT speed captures the number of words post-edited per hour as a measure of productivity. It was obtained by dividing the words edited (length of the respective text) by the time elapsed (task time) and then converting the result to per-hour values. Table 5 shows the summary statistics of the PEMT speed.

Condition	Text	Variable	n	mean	sd
nos	t1	PEMT speed	20	940.853	282.471
nos	t2	PEMT speed	10	1201.389	369.765
nos	t3	PEMT speed	11	853.898	220.689
s	t2	PEMT speed	11	729.944	188.497
s	t3	PEMT speed	10	1040.567	393.146
s	t4	PEMT speed	21	861.932	286.407

Table 5: Summary statistics of the *PEMT speed* with mean and standard deviation (sd).

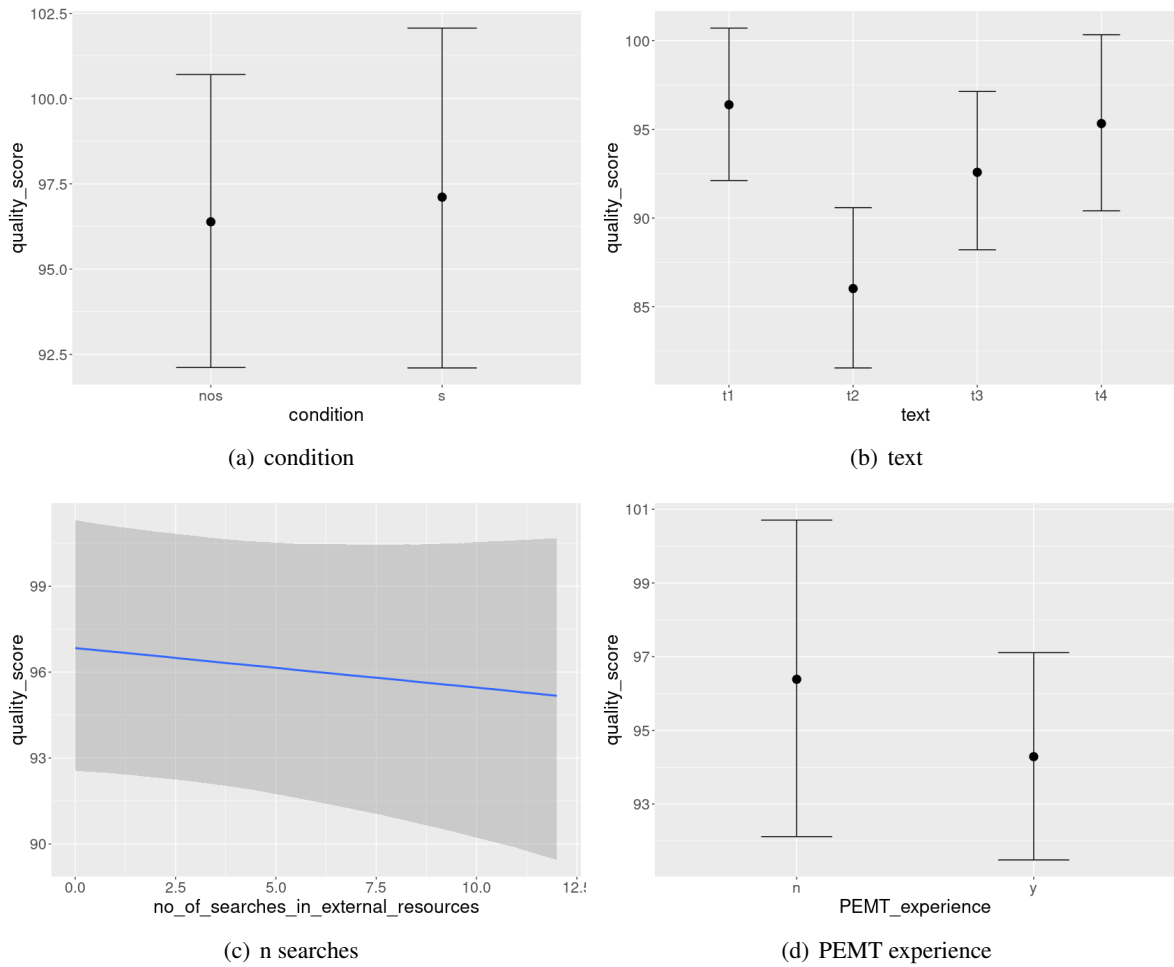


Figure 1: Conditional effects of a) condition, b) text, c) n searches, and d) PEMT experience predictors on *quality score*.

Population-Level Effects			
Predictors	Estimate	CI (95%)	ROPE ↓
Intercept	776.34	[514.94, 1031.95]	0.00%
condition [s]	-137.04	[-201.28, -66.16]	0.00%
text [t2]	88.21	[16.00, 159.16]	3.74%
text [t3]	64.27	[-5.62, 143.43]	15.83%
text [t4]	71.97	[-15.46, 159.45]	15.73%
n searches	-13.01	[-24.74, -1.42]	100%
PEMT experience [y]	254.74	[-34.69, 547.15]	3.92%
Group-Level Effects			
	sd	CI (95%)	
Intercept	287.35	[205.61, 404.32]	
condition [s]	53.27	[2.36, 129.99]	

Table 6: Summary of the fitted model for the *PEMT speed*. ROPE size ± 31.39 .

Table 6 shows the model summary for PEMT speed with a substantial negative effect of the sound condition on the PEMT speed. There are differences across the 4 texts, with a substantial effect

observed for t2. The PEMT experience has a substantial positive effect on productivity.

Figure 2 shows the conditional effects for the PEMT speed. The sound condition decreases PEMT speed in a), there is a large difference across texts in b) with the highest in t2, and an increase in the number of searches decreases the PEMT speed with high uncertainty, in c). As shown in Figure 2 d) having PEMT experience (y) increases productivity, where the difference from no experience (n) has low uncertainty. Figure 5 b) (Appendix) shows the fitted curve with the data points across texts from the productivity model. For t2 the sound condition decreases the PEMT speed, but with t3 there is an increase in speed.

3.3 Cognitive Effort

We define outcome variables for the cognitive effort with the following eye-tracking measures: MFD-ST and MFD-TT. These measures are used as a secondary indicator of the cognitive resources ex-

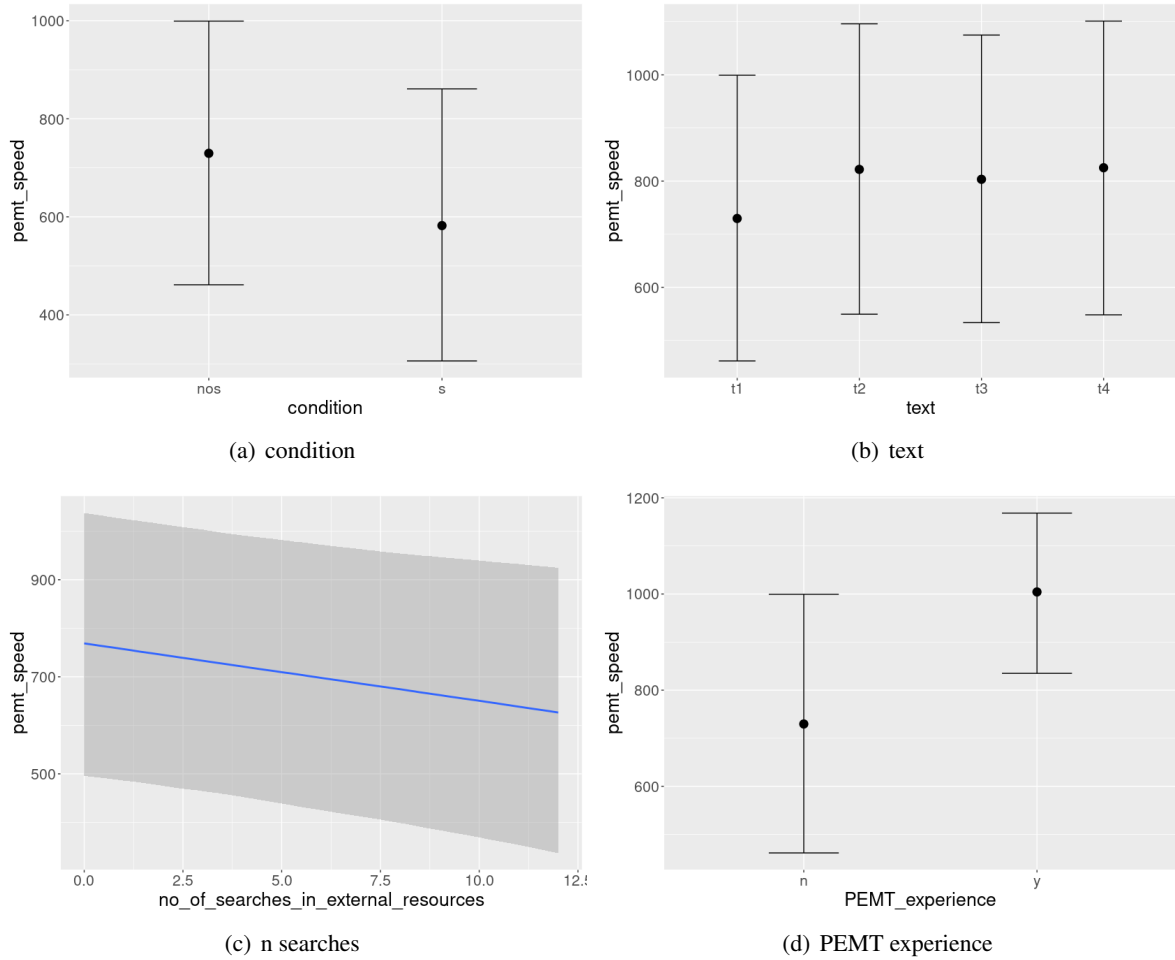


Figure 2: Conditional effects of a) condition, b) text, c) n searches, and d) PEMT experience predictors on *PEMT speed*.

pended by participants, based on the eye-mind assumption (Just and Carpenter, 1980). Mean fixation duration is defined as the total time spent in fixations (keeping the eye stable above a point of focus), divided by the total number of fixations, and is therefore an indication of how long elements of the source and target text were fixated on average. Longer fixations are assumed to indicate higher cognitive effort (Holmqvist and Andersson, 2017). When using a method based on *visual* allocation of attention in an experiment including a listening component, it is important to note that MFD does not reflect how much time the participants spend looking at the screen, which could be assumed to be lower when adding speech synthesis to the process. Rather, MFD reflects how long fixations last on average and is therefore indicative of how effortful processing the text was for participants when they were reading it. Table 7 shows the summary statistics of the MFD-ST.

Table 8 shows the model summary for the MFD-

Condition	Text	Variable	n	mean	sd
nos	t1	MFD_ST	19	298.216	51.833
nos	t2	MFD_ST	9	319.582	61.672
nos	t3	MFD_ST	10	308.829	57.052
s	t2	MFD_ST	10	338.647	56.24
s	t3	MFD_ST	9	315.406	57.277
s	t4	MFD_ST	19	352.568	69.19

Table 7: Summary statistics of the *MFD-ST* with mean and standard deviation (sd).

ST. The sound condition has a non-substantial positive effect on the MFD-ST. There are differences across the texts, with t2 and t4 having the highest effect on the MFD-ST.

Figure 3 shows the conditional effects for the MFD-ST. The sound condition increases the MFD-ST in a) with high uncertainty, there is no large difference across texts in b), the number of searches increases the MFD-ST with high uncertainty in c), and having PEMT experience (y) increases the MFD-ST with low uncertainty.

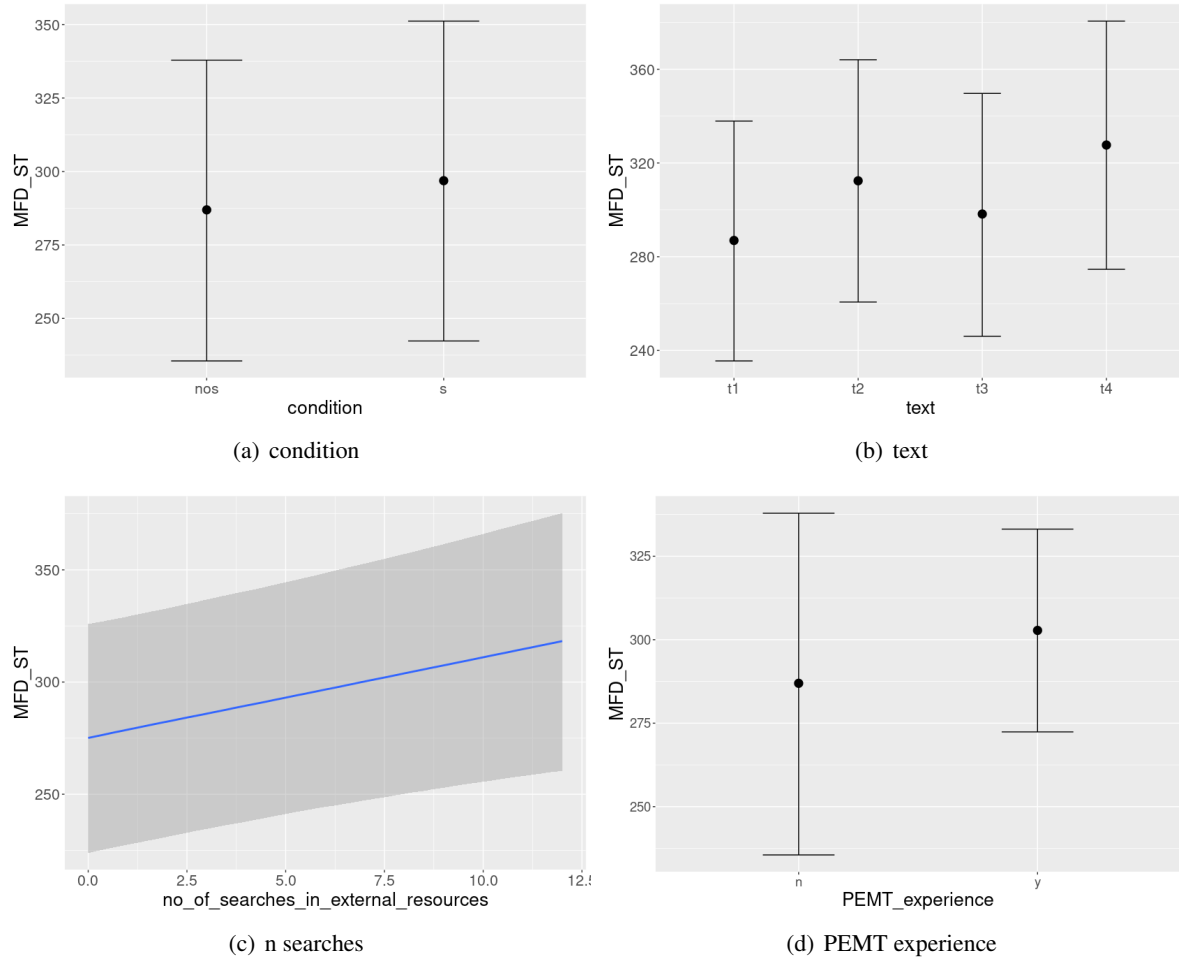


Figure 3: Conditional effects of a) condition, b) text, c) n searches, and d) PEMT experience predictors on *MFD-ST*.

Population-Level Effects			
Predictors	Estimate	CI (95%)	ROPE ↓
Intercept	275.10	[223.84, 325.84]	0.00%
condition [s]	9.85	[-7.52, 27.04]	31.61%
text [t2]	25.49	[9.34, 41.36]	0.00%
text [t3]	11.24	[-4.63, 26.97]	25.37%
text [t4]	40.66	[20.76, 60.30]	0.00%
n searches	3.61	[0.94, 6.22]	99.63%
PEMT experience [y]	16.08	[-42.96, 75.15]	14.78%
Group-Level Effects			
	sd	CI (95%)	
Intercept	55.13	[38.11, 79.49]	
condition [s]	20.45	[3.22, 38.36]	

Table 8: Summary of the fitted model for the *MFD-ST*. ROPE size ± 6.14 .

Table 9 shows the summary statistics of the MFD-TT. Table 10 shows the model summary for the MFD-TT. The sound condition has a substantial negative effect on the MFD-TT. There are substan-

tial differences across the texts, with t4 having the highest effect on the MFD-TT.

Condition	Text	Variable	n	mean	sd
nos	t1	MFD_TT	19	382.189	62.845
nos	t2	MFD_TT	9	416.568	69.55
nos	t3	MFD_TT	10	413.299	69.448
s	t2	MFD_TT	10	415.418	77.357
s	t3	MFD_TT	9	378.956	63.564
s	t4	MFD_TT	19	421.6	76.602

Table 9: Summary statistics of the *MFD-TT* with mean and standard deviation (sd).

Figure 4 shows the conditional effects for the MFD-TT. The sound condition decreases the MFD-TT in a) with high uncertainty, there is no large difference across texts in b), the number of searches is associated with a small increase in MFD-TT with high uncertainty in c), and having PEMT experience (y) decreases the MFD-TT in d) with low uncertainty but a large overlap with *no experience* (n). Figure 5 (Appendix) shows the fitted curve with the

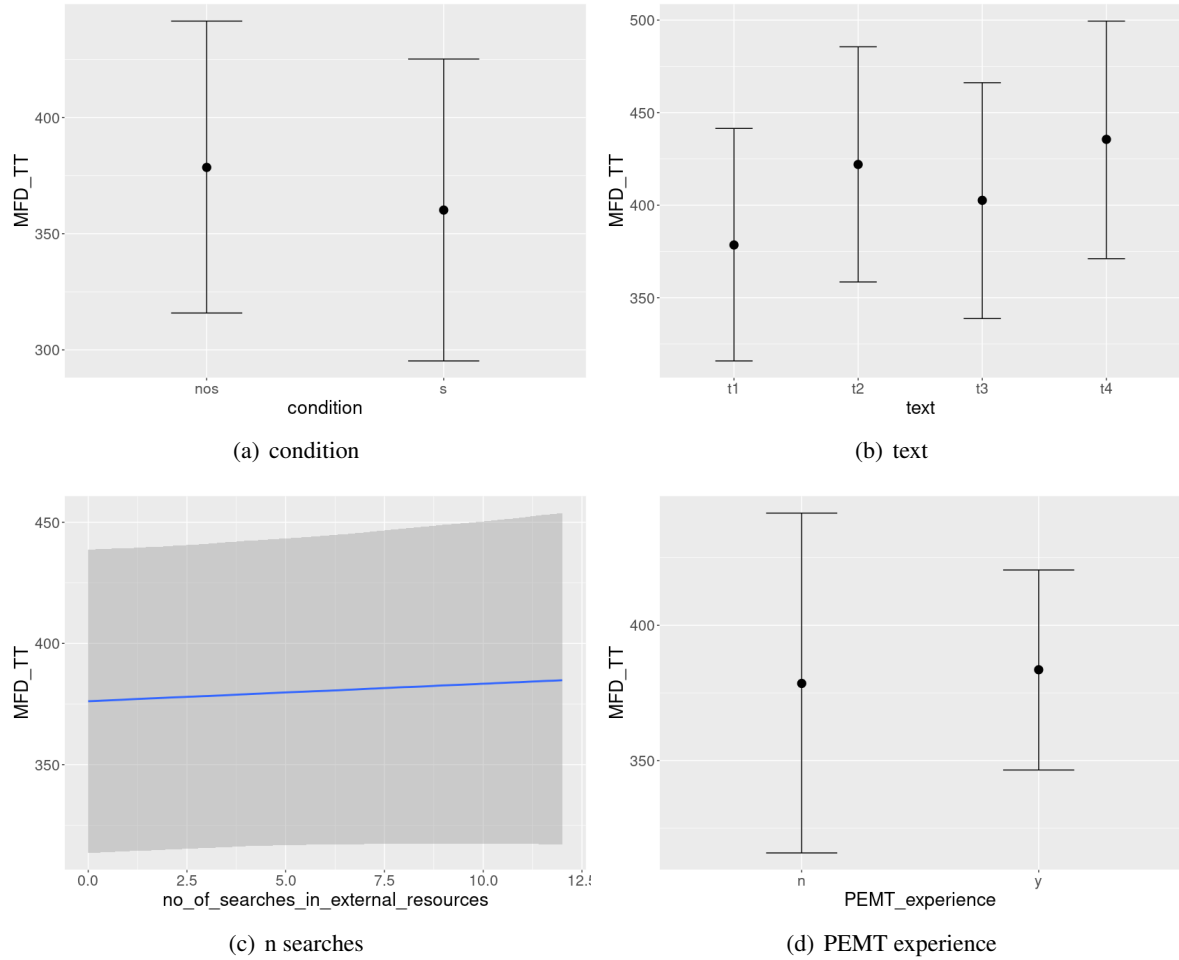


Figure 4: Conditional effects of a) condition, b) text, c) n searches, and d) PEMT experience predictors on *MFD-TT*.

Population-Level Effects			
Predictors	Estimate	CI (95%)	ROPE ↓
Intercept	376.07	[313.86, 438.67]	0.00%
condition [s]	-18.34	[-32.42, -4.19]	3.41%
text [t2]	43.28	[28.45, 58.14]	0.00%
text [t3]	23.96	[9.32, 38.68]	0.00%
text [t4]	56.95	[38.70, 75.45]	0.00%
n searches	0.73	[-1.74, 3.17]	100%
PEMT experience [y]	5.03	[-69.48, 79.38]	16.25%
Group-Level Effects			
	sd	CI (95%)	
Intercept	68.22	[48.20, 97.29]	
condition [s]	9.74	[0.61, 22.42]	

Table 10: Summary of the fitted model for the *MFD-TT*. ROPE size ± 7.02 .

data points across texts from c) *MFD-ST*, and d) *MFD-TT*. Figure c) shows that the sound condition increases the *MFD-ST* for t2, but decreases it for t3. The same pattern is observed for *MFD-TT* in

d), where the sound condition is associated with an increase for t2, and a decrease for t3.

4 Discussion

The results of our experiment on using speech synthesis for PEMT indicate that (1) differences in quality between conditions were small; (2) participants were slower when using speech synthesis; and (3) participants expended less cognitive effort in TT when using speech synthesis, as reflected in their fixation data. More specifically, the presence of speech had a substantial negative effect on the *MFD-TT*, meaning that overall the cognitive effort spent by translators reading the target text was reduced. This may mean that hearing the target text was considered by translators to be a reliable source of information when checking PEMT. We report a non-substantial positive effect on the *MFD-ST* variable, indicating that the processing of the source text does not change much and only increases slightly. We do not believe this to be due

to a lack of trust in the speech synthesis, given the results for the MFD-TT, but that speech use may be more worthwhile in the TT. This is also suggested by the answer to the perception questionnaire we distributed at the end of the experiment (Ciobanu et al., forthcoming) where the most reported on advantages of using speech were improved style (11/21) and error detection (9/21). It may also be that listening to the TT causes the translators to expend more cognitive effort on the ST, but this would require a separate analysis. The decrease in productivity might reflect the fact that listening to the text is an additional step to be carried out in the workflow. Moreover, as all participants but one were first-time users of speech synthesis in PEMT, productivity losses can reasonably be expected to diminish as users become more familiar with the tool. A longitudinal study would surely provide useful data in this regard. Related to this but apart from the effect of the sound condition, we also found that PEMT experience has a substantial positive effect on PEMT speed, indicating that translators with previous PEMT experience work faster than those without. The effect of the number of searches is non-substantial for all outcome variables. We recorded no substantial change in quality, but there is a perceived improvement in style and error detection for some of the participants as reported in (Ciobanu et al., forthcoming). The loss in productivity may be reduced following longer exposure to speech synthesis. This, coupled with the substantial decrease in cognitive effort in the TT, point to a potential support that a speech-enabled mode of working can offer translators.

5 Conclusions and Future Work

We quantified the impact of text-to-speech on PEMT for the English-German language pair. We introduce a Bayesian hierarchical model to tackle issues with data scarcity. The introduction of the sound condition on the PEMT workflow has a non-substantial positive effect on quality, a substantial negative effect on PEMT speed, and a non-substantial positive effect on MFD-ST and substantial negative effect on the MFD-TT for cognitive effort. The effect of the number of searches is non-substantial for all outcome variables. The text ID together with the sound condition has an effect on all of the measurements, which may be explained by the standard measurements of text complexity we used, which do not take into account semantics

and might not sufficiently reflect textual differences, especially regarding translation difficulty.

For future work, we will measure the relation between text complexity evaluated with newer readability formulas based on fine-grained linguistic features and translation quality/productivity (Dai and Liu, 2024), investigate in more detail the relation between translation experience and translation quality/productivity, the relation between productivity and the number of searches performed, and quantify the observable changes in individual PEMT workflows created by our participants' access to speech synthesis.

Acknowledgements

This project received funding from the Imminent Research Grants programme. We also thank the project participants and the reviewers for their feedback.

References

- Burchardt, Aljoscha. 2013. Multidimensional quality metrics: a flexible system for assessing translation quality. In *Proceedings of Translating and the Computer 35*, London, UK, November 28-29. Aslib.
- Bürkner, Paul-Christian. 2017. brms: An R package for Bayesian multilevel models using Stan. *Journal of Statistical Software*, 80(1):1–28.
- Ciobanu, Dragoş, Valentina Ragni, and Alina Secară. 2019. Speech Synthesis in the Translation Revision Process: Evidence from Error Analysis, Questionnaire, and Eye-Tracking. *Informatics*, 6(4)(51), December.
- Ciobanu, Dragoş, Miguel Rios, Alina Secară, Justus Brockmann, Raluca-Maria Chereji, and Claudia Wiesinger. forthcoming. The impact of speech synthesis on cognitive effort, productivity, quality, and perceptions during post-editing machine translation (PEMT). *Revista Tradumática: translation technologies*.
- Dai, Guangrong and Siqi Liu. 2024. Towards predicting post-editing effort with source text readability: An investigation for english-chinese machine translation. *The Journal of Specialised Translation*, (41):206–229, Jan.
- ELIA, EMT, EUATC, FIT EUROPE, GALA, LIND, and Women in Localization. 2023. 2023 European Language Industry Survey. Trends, expectations and concerns of the European language industry. Technical report.
- Gelman, Andrew and Jennifer Hill. 2007. *Data analysis using regression and multilevel/hierarchical models*,

- volume Analytical methods for social research. Cambridge University Press, New York.
- Gelman, Andrew, John B. Carlin, Hal S. Stern, and Donald B. Rubin. 2004. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2nd ed. edition.
- Gius, Evelyn, Jan Christoph Meister, Malte Meister, Marco Petris, Mareike Schumacher, and Dominik Gerstorfer. 2023. *CATMA 7 (Version 7.0)*. Zenodo.
- Guerberof-Arenas, Ana and Antonio Toral. 2022. Creativity in translation: machine translation as a constraint for literary texts. *Translation Spaces*, 11(2):184–212, November. Publisher: John Benjamins Publishers.
- Herbig, Nico, Tim Düwel, Santanu Pal, Kalliopi Meladaki, Mahsa Monshizadeh, Antonio Krüger, and Josef van Genabith. 2020. MMPE: A Multi-Modal Interface for Post-Editing Machine Translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1691–1702, Online, July. Association for Computational Linguistics.
- Holmqvist, Kenneth and Richard Andersson. 2017. *Eye tracking: a comprehensive guide to methods, paradigms and measures*. Lund Eye-Tracking Research Institute, Lund, Sweden, 2nd edition edition.
- Just, Marcel A. and Patricia A. Carpenter. 1980. A theory of reading: From eye fixations to comprehension. *Psychological Review*, 87(4):329–354, July. Publisher: American Psychological Association.
- Kenny, Dorothy. 2020. Machine Translation. In Baker, Mona and Gabriela Saldanha, editors, *Routledge Encyclopedia of Translation Studies*, Routledge Handbooks in Translation and Interpreting Studies, pages 305–310. Routledge, 3 edition.
- Kim, Yu Yeon, Aluko Ademola, Jeong Hyeun Ko, and Hee Sook Kim. 2022. Knuir at the ntcir-16 rcir: Predicting comprehension level using regression models based on eye-tracking metadata.
- Krings, Hans P. 2001. *Repairing texts: empirical investigations of machine translation post-editing processes*. The Kent State University Press, Ohio.
- Kruschke, John K. 2018. Rejecting or accepting parameter values in bayesian estimation. *Advances in Methods and Practices in Psychological Science*, 1(2):270–280.
- Läubli, Samuel, Chantal Amrhein, Patrick Düggin, Beatriz Gonzalez, Alena Zwahlen, and Martin Volk. 2019. Post-editing Productivity with Neural Machine Translation: An Empirical Assessment of Speed and Quality in the Banking and Finance Domain. *arXiv:1906.01685 [cs]*, June. arXiv: 1906.01685.
- McElreath, Richard. 2016. *Statistical rethinking: a Bayesian course with examples in R and Stan*. Number 122 in Chapman & Hall/CRC texts in statistical science series. CRC Press/Taylor & Francis Group, Boca Raton. largely / videos.
- Moorkens, Joss and Sharon O’Brien. 2017. Assessing user interface needs of post-editors of machine translation. In Kenny, Dorothy, editor, *Human Issues in Translation Technology*, pages 109–130. Routledge, London.
- Moorkens, Joss, 2018. *Eye tracking as a measure of cognitive effort for post-editing of machine translation*, page 55–70. John Benjamins Publishing Company, September.
- O’Brien, Sharon. 2009. Eye-tracking in translation process research: Methodological challenges and solutions. In Mees, Inger M., Susanne Göpferich, and Fabio Alves, editors, *Methodology, Technology and Innovation in Translation Process Research: A Tribute to Arnt Lykke Jakobsen*, pages 251–266. Samfundslitteratur.
- Plitt, Mirko and François Masselot. 2010. A Productivity Test of Statistical Machine Translation Post-Editing in a Typical Localisation Context. *The Prague Bulletin of Mathematical Linguistics*, 93:7. Num Pages: 7 Place: Prague, Poland Publisher: De Gruyter Poland.
- Silva, Breno B., David Orrego-Carmona, and Agnieszka Szarkowska. 2022. Using linear mixed models to analyze data from eye-tracking research on subtitling. *Translation Spaces*, June. © John Benjamins Publishing Company.
- Taylor, Paul. 2009. *Text-to-Speech Synthesis*. Cambridge University Press.
- Wiesinger, Claudia, Justus Brockmann, Alina Secară, and Dragoş Ciobanu. 2022. Speech-enabled machine translation post-editing in the context of translator training. In Kornacki, Michał and Gary Massey, editors, *Contextuality in Translation and Interpreting. Selected Papers from the Łódź-ZHAW Duo Colloquium on Translation and Meaning 2020–2021*, volume 70 of *Łódź Studies in Language*. Peter Lang.

A Model Formulas

In this section, we show the brms formulas for each outcome variable Y .

Quality score outcome: quality score, first level predictors: condition, text, n searches, PEMT experience, and second level predictors: condition. brms formula:

$$\text{quality_score} \sim 1 + \text{condition} + \text{text} + \text{n_searches} + \text{pemt_experience} + (1 + \text{condition} | \text{participant})$$

PEMT speed productivity outcome: PEMT speed, first level predictors: condition, text, n searches, PEMT experience, and second level predictors: condition. brms formula:

$$\text{pemt_speed} \sim 1 + \text{condition} + \text{text} + \text{n_searches} + \text{pemt_experience} + (1 + \text{condition} \mid \text{participant})$$

MFD-ST outcome: MFD-ST, first level predictors: condition, text, n searches, PEMT experience, and second level predictors: condition. brms formula:

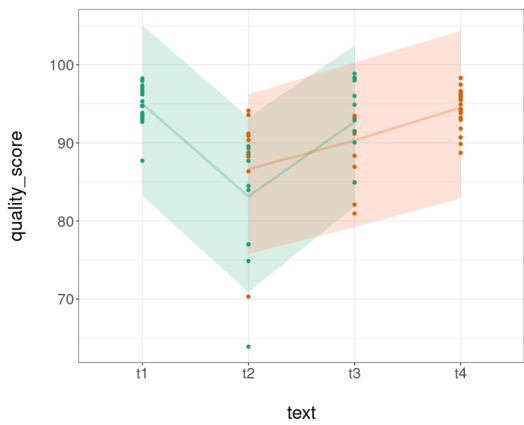
$$\text{MFD_ST} \sim 1 + \text{condition} + \text{text} + \text{n_searches} + \text{pemt_experience} + (1 + \text{condition} \mid \text{participant})$$

MFD-TT outcome: MFD-TT, first level predictors: condition, text, n searches, PEMT experience, and second level predictors: condition. brms formula:

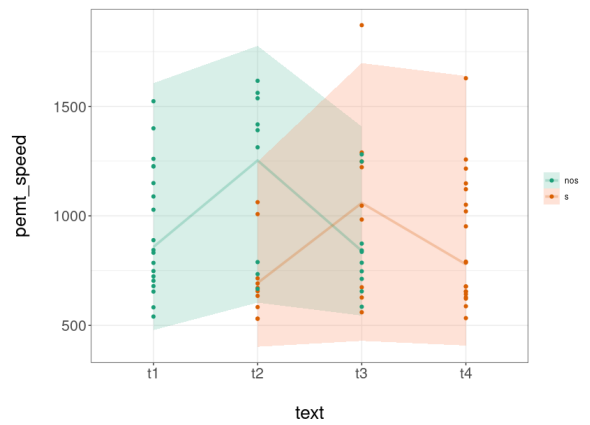
$$\text{MFD_TT} \sim 1 + \text{condition} + \text{text} + \text{n_searches} + \text{pemt_experience} + (1 + \text{condition} \mid \text{participant})$$

B Fitted Models

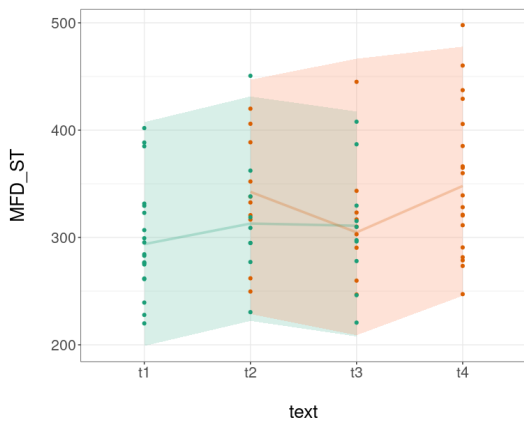
C ROPE for the Sound Condition



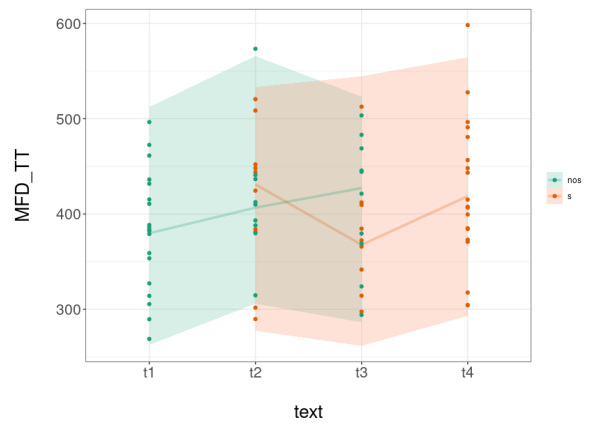
(a) Quality



(b) PEMT speed



(c) MFD-ST



(d) MFD-TT

Figure 5: Fitted models across texts on each condition for: a) *Quality*, b) *PEMT speed*, c) *MFD-ST*, and d) *MFD-TT*. Fit curve with posterior predictions from the model, uncertainty bars with 95% CI, and data points.

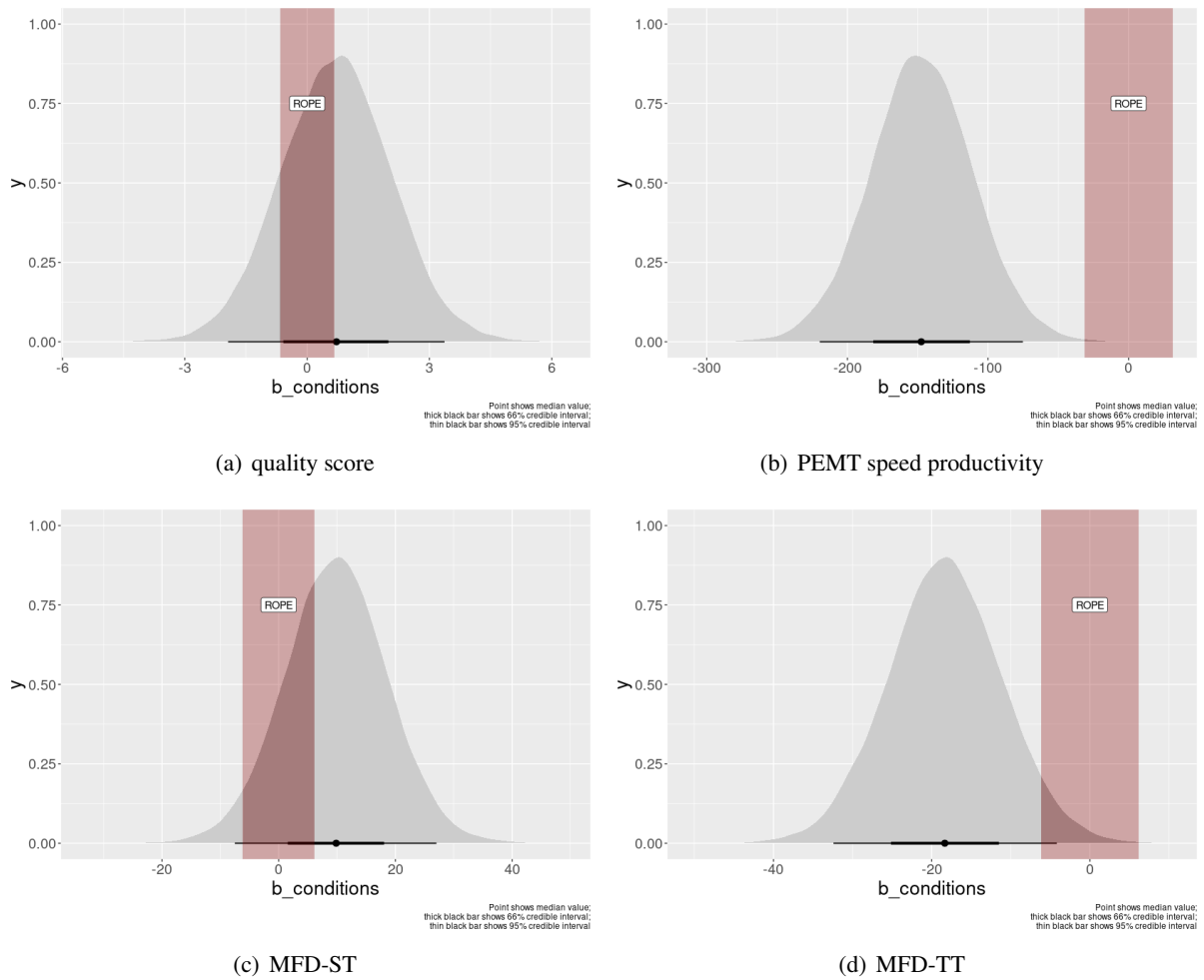


Figure 6: Proportion of the sound condition effect that falls into the ROPE for each outcome variable: a) *quality*, b) *PEMT speed*, c) *MFD-ST*, and d) *MFD-TT*. Point median value, thin bar 95% CI, and thick bar 66%CI.

Evaluation of intralingual machine translation for health communication

Silvana Deilen¹, Ekaterina Lapshinova-Koltunski¹, Sergio Hernández Garrido¹,
Christiane Maaß¹, Julian Hörner², Vanessa Theel³, Sophie Ziemer⁴

¹ University of Hildesheim, ² Wort & Bild Verlag, ³ SUMM AI, ⁴ Johannes Gutenberg University Mainz
¹deilen, lapshinovakoltun, hernandezs, maass@uni-hildesheim.de,
²j.hoerner@wubv.de, ³vanessa@summ-ai.com, ⁴sziemer@students.uni-mainz.de

Abstract

In this paper, we describe results of a study on evaluation of intralingual machine translation. The study focuses on machine translations of medical texts into Plain German. The automatically simplified texts were compared with manually simplified texts (i.e., simplified by human experts) as well as with the underlying, unsimplified source texts. We analyse the quality of outputs from three models based on different criteria, such as correctness, readability, and syntactic complexity. We compare the outputs of the three models under analysis between each other, as well as with the existing human translations. The study revealed that system performance depends on the evaluation criteria used and that only one of the three models showed strong similarities to the human translations. Furthermore, we identified various types of errors in all three models. These included not only grammatical mistakes and misspellings, but also incorrect explanations of technical terms and false statements, which in turn led to serious content-related mistakes.

1 Introduction

In Germany, according to recent studies in the field of Public Health, over half of the population reports having difficulties with health-related topics (Schaeffer et al., 2021). For that reason, the promotion of health literacy (knowledge and compe-

tences to access, understanding, appraise and apply medical information) has turned into an important task for the German health system (Schaeffer et al., 2018) (for an extensive definition of health literacy, see Sørensen et al. (2012)). In this context, recent research has underlined the need for accessible communication in the medical domain to effectively promote health literacy and consequently assist patients navigating the health system and improve patient understanding, engagement and compliance with medical recommendations (Ahrens et al., 2022; Blechschmidt, 2021; Schaeffer et al., 2021). Plain German is a prominent form of accessible communication that has gained relevance in health communication scenarios (Schaeffer et al., 2018).

Although there is an urgent need for translations into Plain German, there is also a gap in qualified and experienced human translators (Maaß, 2020). Moreover, there is a lack in computer-aided translation (CAT) tools and machine translation systems for this kind of intralingual translation. Unfortunately, little is known about existing systems and their performance for different texts that are required to be translated into Plain German, as for instance, texts in health communication that we focus on.

In our study, we evaluate machine translations of medical texts into Plain Language. The source texts, as well as reference human translations, are derived from the website of the German health magazine *Apotheken Umschau*. We analyse machine-translated texts produced with three models comparing them with human translations from the magazine’s website. Besides that, we compare all translations with the underlying sources.

In the following, we present the results of the qualitative and quantitative analysis. Section 2 de-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

scribes related work. In Section 3, we present our research design including corpus and the methods used. Section 3.3 presents the results of our analyses, while we discuss those results, as well as limitations and possible extensions in the future in Section 4.

2 Related Work

2.1 Plain Language

Both Easy Language and Plain Language are complexity-reduced language varieties which aim to improve readability and comprehensibility of texts (Bredel and Maaß, 2016; Maaß, 2020). They are used in different communication scenarios, e.g. in legal communication (Maaß and Rink, 2021) or health communication (Ahrens et al., 2022), and have different target groups (Maaß and Schwengber, 2022). While Easy Language is characterized by a maximally reduced complexity on all language levels and is mainly intended for people with communication impairments and disabilities, the grammatical and textual features of Plain Language are closer to the standard language and are mainly a means to open expert contexts for lay people (Maaß, 2020). Therefore, the main target audience of Plain Language are non-experts with average or slightly below average language or reading skills (Maaß, 2020). In Germany, Easy Language has become a subject of scientific research since 2014 with rapidly growing output of publications in the following years (Maaß et al., 2021; Deilen et al., 2023b). The studies point in two basic directions: studies on text qualities and possible barriers in various forms of communication on the one side (Rink, 2019) and studies on comprehensibility and recall by different target groups on the other (Gutermuth, 2020; Deilen, 2021).

Unlike Easy Language, Plain Language is a dynamic variety. Plain Language does not have a fixed set of rules, but the linguistic complexity of Plain Language texts is adapted to the needs of the intended audience in a specific target situation (Bredel and Maaß, 2016; Maaß, 2020). Therefore, Plain Language is a flexible concept that varies depending on the presumed reading skills of its target group (Maaß, 2020). In comparison to Easy Language, Plain Language has the advantage of having less to no stigmatizing features, which is one of the reasons why it is also more acceptable than Easy Language. However, due to the higher degree

of linguistic complexity, Plain Language texts are far less comprehensible than Easy Language texts and therefore not necessarily accessible for people with very low literacy skills. Maaß (2020) therefore models the variety Easy Language Plus, which is situated between Easy Language and Plain Language and strikes a balance between comprehensibility and acceptability¹.

2.2 Accessibility in Medical Domain

In 2016, the Health Literacy Survey (HLS-GER) revealed that more than half of the German population (54.3%) encounters significant challenges in locating, understanding, appraising, and effectively using health-related information. These findings, which according to Schaeffer et al. (2017) were “significantly worse than expected” increased the awareness of the need for accessible health information, which in turn led the development of the National Action Plan Health Literacy (Schaeffer et al., 2018). According to this plan, one approach to increase health literacy in Germany is providing information in Plain Language, i.e., in a complexity-reduced variety of German. With the release of updated data from the second Health Literacy Survey (HLS-GER 2) in 2021, the importance of Plain Language in German health communication has been underscored, as it has shown that even more people (58.8%) encounter difficulties navigating the healthcare system. As a remedy for low health literacy, both practitioners and researchers increasingly advocate for the use of Plain Language. One of the most prominent examples of implementing this approach is the *Apotheken Umschau*. The *Apotheken Umschau*, which is Germany’s leading health publisher and the largest consumer medium in the German-speaking area with a traffic of 6.68 m. visits and 49.11 m. page impressions per month², has so far published more than 220 texts in Plain Language on their website in a co-operation with the Research Centre for Easy Language (University of Hildesheim). By offering information in both standard German and Plain German, their goal is to provide accessible and reliable information on illnesses, medications, and preventive healthcare

¹It should be noted that Plain Language is an international concept and not language-bound. However, in this paper we only focus on Plain Language in Germany, also called Plain German.

²<https://ausweisung-digital.ivw.de/>, retrieved 15.03.2024

with minimal barriers for all individuals.

2.3 NLP for Plain Languages

In Easy and Plain Language translation, which both belong to the domain of intralingual translation (Hansen-Schirra et al., 2020a), the potentials of using CAT tools are still a major research desideratum. There are some studies that have discussed the challenges of using CAT tools in intralingual translation compared to interlingual translation (Hansen-Schirra et al., 2020b; Spring et al., 2023; Kopp et al., 2023). For example, in contrast to interlingual translation, in intralingual translation there is usually no 1:1 correspondence between source and target sentences, which in turn means that the sentence alignment process has to be done or corrected manually by the translator, which increases the workload for translators instead of reducing it.

While there are plenty of studies on automatic text simplification methods that aim to automatically convert a text into another text that is easier to understand, while ideally conveying the same message as the source text, which contributes to textual accessibility (Sheang and Saggion, 2021; Maddela et al., 2021; Martin et al., 2020; Saggion, 2017), most of them do not consider the needs of the target audience. Scarton and Specia (2018) showed that using target audience oriented data helps to build better models for automatic text simplification using the Newsela corpus³. However, this corpus contains news texts only, whereas we are looking into the medical discourse, where texts in Plain Language enable accessibility to health literacy.

Ondov et al. (2022) surveyed the literature in the field of automated methods for biomedical text simplification and found that one major challenge in this field is the lack of high-quality parallel text data, which so far impedes the development of fully automated biomedical text simplification methods.

Specific problems of automatic systems for intralingual translation, e.g. copying source segments into the output, were addressed by Säuberli et al. (2020) and Spring et al. (2023), who showed that pretrained and fine-tuned NMT models have promising results in automatic text simplification. However, as stated by Anschütz et al. (2023), even though there are improvements in the systems

of automated intralingual translation, the outputs might, so far, not be used by the target groups directly. Nevertheless, they may serve as a draft for professional intralingual translators to reduce their workload. Deilen et al. (2023a) drew similar conclusions for the outputs produced with ChatGPT. The authors investigated the feasibility of using ChatGPT for intralingual translation. They analysed the quality of the generated texts according to such criteria as correctness, readability, and syntactic complexity. Their results indicated that the texts produced by ChatGPT were easier than the standard source texts, but the content was not always rendered correctly. Besides that, the automated intralingual output did not fully meet the standards which human translators follow. In the present study, we follow a similar approach. However, while the authors analysed intralingual translation into German Easy Language, a maximally simplified and strictly controlled language variety adapted to the needs of people with reading impairments, we focus on translation into Plain German (see 2.1). Besides that, we focus on medical texts, whereas the authors translated citizen-oriented administrative texts. Moreover, we investigate the feasibility of a tool which was specifically trained for intralingual translation into Easy and Plain Language instead of using a chatbot designed for various tasks.

3 Research Design

3.1 Data Collection

Our dataset contains 200 parallel texts selected from the website of the German health magazine *Apotheken Umschau*⁴. The texts cover a broad range of topics, such as breast cancer, vaccination, long COVID, food poisoning, first aid and others. For all texts in the sample, a human translation into Plain Language was already available. Both the source texts and the human translations were reviewed by medical or pharmaceutical professionals from the editorial team of *Apotheken Umschau* and comply with the guidelines of evidence-based medicine. Content accuracy is therefore guaranteed for the sample. Furthermore, the human translations also comply with a practical concept for Plain Language for this specific health information scenario, which was established by the Research Centre for Easy Language⁵. We split the data into

³<https://newsela.com/data>

⁴apotheken-umschau.de

⁵www.uni-hildesheim.de/leichtesprache

test and train sets: 30 texts were selected to serve as test data in our evaluation study and the remaining 170 texts were used as training data for two of the three tested systems. The sample of 30 texts was translated using the machine translation system SUMM AI⁶. At the time of the study, SUMM AI was the only tool known to us for intralingual translation from standard German into Easy and Plain German and the only one with a specific focus on health communication texts.

In our study, we compared three different models of SUMM AI: the baseline model and two further fine-tuned models. The baseline model of our study was the already existing beta-model for Plain Language provided by SUMM AI. The model is a fine-tuned large language model (LLM) that was trained with in-house data and further rule-based approaches. In comparison to the baseline model, two further models (model 1 and model 2 hereafter) were enriched by SUMM AI with the training data (170 parallel texts out of 200 selected). The data was aligned and adapted according to the practical concept of the Research Centre for Easy Language. While the baseline model and model 1 have the same underlying LLM, model 2 is distinguished by a different underlying LLM.

We investigate which of the three models yield better results in translating standard German into Plain German. For this, the 30 texts from the test set are translated with the three models under analysis⁷.

3.2 Data Analysis

3.2.1 Analysis Steps

The resulting machine translations (three per each texts) are also compared with the existing human translations and the underlying sources. We follow the evaluation criteria suggested by Deilen et al. (2023a), which is one of the few studies known to us that evaluates intralingual machine translation. We assess machine translations for the correctness of the content, the readability of the texts, and their syntactic complexity. Readability as well as syntactic complexity are also assessed for human translations and source texts. We then

⁶SUMM AI is a tool for translating texts into Easy German and Plain German. The company SUMM AI offers different licenses for freelancers, authorities and companies, see <https://summ-ai.com> for more details.

⁷The whole dataset will be published on GitHub. The GitHub repository will contain the selected texts (sources, human and machine translations), including the raw data, the parsed data (conllu) and the Textlab analyses per text.

compare sources, human, and machine translations according to these two criteria.

3.2.2 Correctness

The content of the machine-generated texts was first analysed for correctness. This content evaluation was done manually, whereby each text was assessed independently by two researchers, who checked whether the medical information in the target text is still valid despite reduction of complexity and shortening of information. In cases where an accurate assessment required specialized knowledge, a healthcare professional from the *Apotheken Umschau* team was consulted. No quantitative error analysis was performed. Consequently, a translation was already considered incorrect if it contained one content-related error. This is because the study seeks insights into who artificial intelligence (AI) powered translation tools are suitable for: translators, content providers, or Plain German end users (for an overview over end users, see Bredel/Maaß, 2016 and Maaß 2020). In order for machine translation into Easy or Plain Language to be safely usable by end users, the target texts must not contain errors. The presence of errors in the target texts therefore indicates usability for users other than the end users.

3.2.3 Readability

We also compared the readability of human and the machine translations, as well as of the source texts. For this, we use the Hohenheim Comprehensibility Index (HIX). The HIX is a meta index that calculates the readability of a text taking into account the four major readability formulas common in Easy Language Research (Bredel and Maaß, 2016, p. 61ff). They include the Amstad index, the simple measure of gobbledygook (G-SMOG) index, the Vienna non-fictional text formula (W-STX) and the readability index (LIX), with a HIX of 0 indicating extremely low comprehensibility and a HIX of 20 extremely high comprehensibility (for further details see: <https://klartext.uni-hohenheim.de/hix>). The benchmark for a text to be classified as a text in Easy German, which is the least complex variety of German, is set at 18 points (Rink, 2019). As Plain German is more complex than Easy German, we suggest setting the benchmark for Plain German at 16 points.

3.2.4 Syntactic Complexity

We operationalised syntactic complexity as a distribution of specific syntactic relations, i.e. specific clauses. We automatically identified syntactic relations using dependency parsing that we obtained with the Stanford NLP Python Library Stanza (v1.2.1)⁸ with all the models pre-trained on the Universal Dependencies v2.5 datasets. Our list of selected structural categories include adnominal clauses or clausal modifiers of noun (acl), adverbial clause modifiers (advcl), clausal components (ccomp), clausal subjects (csubj), open clausal elements (xcomp) and parataxis relation (parataxis). These selected categories are all listed under the clause dependents⁹ in the Universal Dependency. More details on dependency relations and their definitions across languages can be found in De Marneffe et al. (2021). We collected and compared the distribution frequencies of these categories in the three subcorpora under analysis (source texts, human translations, and machine translations). We interpreted the results based on the assumption that the higher the number of these dependency relations in the corpus, the more complex the texts contained in these sub-corpora are.

3.2.5 Automatic Evaluation Measures

We also applied SARI (Xu et al., 2016), which is a quantitative measure to evaluate automatic text simplification systems. The metric “compares system output against references and against the input sentence” (Xu et al., 2016) and is normally used for evaluation of automatic text simplification models but could also be used to evaluate intralingual translation.

While SARI is normally calculated on a sentence basis, this is not possible in the case of Plain Language since there usually is no sentence-to-sentence alignment but rather an alignment on paragraph level. To calculate these metrics, we aligned the source texts, machine translations, and human translations on a paragraph level and assessed their alignment quality. Since the translation into Plain Language compared to interlingual translation is significantly more liberal in terms of which information is translated, adequate alignment was difficult and only possible for 263 of 946 segments.

⁸<https://stanfordnlp.github.io/stanza/index.html>

⁹<https://universaldependencies.org/u/dep/>

3.2.6 Translation Comparison

In the last step, we compared the performance of the systems taken all criteria together. After that, we compared them with the existing human translations, as well as the underlying source texts. For this, we used an explorative multivariate technique called Correspondence Analysis (CA) performed with the package `ca` in R environment (R Core Team, 2017, R version 3.6.1).

Correspondence analysis (Greenacre, 2007) helps to explore relations between variables in a data set (both those constituting the rows and those in columns) and summarises and visualises data in a two-dimensional plot. We use CA to see which variables, in our case subcorpora representing source texts (source), human translations (human) and the three machine-translated outputs (baseline, model 1 and model 2), have similarities and how these subcorpora correlate with the analysed features (HIX values, syntactic structures) contributing to the similarities. Weighted Euclidean distances, termed the χ^2 distances are measured on the basis of the distributions of these feature across the five subcorpora under analysis. The row (subcorpora) and the column (features) projections are then plotted on the same graph. The larger the differences between the subcorpora, the further apart they are on the map. Proximity between subcorpora and features in the merged map is an approximation of the correlation between them. The position of the dots (subcorpora) and triangles (features) indicates the relative importance of a feature for a subcorpus (see Figure 4). With the help of this technique, we will observe which texts are more similar between each other.

3.3 Results

3.3.1 Correctness

The analysis of the correctness of the machine translations showed that from the baseline model, only one of the 30 texts was correctly translated. The other 29 texts showed problems with regard to their correctness in different aspects. Model 1 yielded similar results, with only two out of thirty texts being classified as correct. For model 2, however, we found that 15 out of 30 texts were translated correctly. Overall, the results are disparate and inconsistent. The texts do not follow a uniform structure and are not action oriented. In practice, they would have to be completely

post-edited. We encounter grammatical errors and misspellings, omissions of relevant pronouns or words, incorrect explanations of technical terms, incorrect statements and advice, wrong segmentation of compounds, etc. It should be emphasized, once again, that no quantitative evaluation was performed because the mere presence of the errors themselves was considered a risk for the primary users. Furthermore, so far, we have not classified or ranked the error types based on severity levels, but we plan to do so in our future work (see Section 4). Some examples of the errors we found are given in the following.

- Missing segmentation signs: In some cases, segmentation signs would facilitate the processing, but the tool fails to apply them. This is especially true for polymorphemic compounds (i.e., compounds consisting of at least three free morphemes), such as *"Nasennebenhöhlenentzündungen"* (model 2) (*inflammation of the sinus cavities*), in which indicating the morpheme boundaries would have reduced the compound's complexity.
- Redundancies and unreasonable statements:
 - *"Bei Männern kann eine Blasenentzündung auch die Prostata entzünden. Oder die Prostata entzündet sich. Dann kann sich die Prostata entzünden."* (model 1) (*In men, a bladder infection can also inflame the prostate. Or the prostate becomes inflamed. Then, the prostate becomes inflamed.*)
 - *"Dann kann eine Person eine Nierenbeckenentzündung oder eine Nierenbeckenentzündung bekommen."* (model 1) (*Then, a person can get an urinary tract infection or an urinary tract infection.*)
 - *"Eine Insekten-Stich ist eine allergische Reaktion auf einen Insekten-Stich."* (model 1) (*An insect bite is an allergic reaction to an insect bite.*)
 - *"Frauen haben oft eine Blasenentzündung, weil sie oft auf Toilette müssen."* (model 1) (*women often have a bladder infection because they often have to go to the toilet.*)

- Lexico-semantic errors:

- *"Viele Menschen nehmen zu wenig Schlaf"* (model 1) (*Many people take too little sleep.*)
- *"Wenn andere Menschen sich Sorgen um Sie machen, ist auch ein Zeichen."* (baseline) (*When other people are worried about you, is also a sign.*)

- Omission of reflexive pronouns:

- *"Dann können sie gut konzentrieren"* (model 1) (*Then they can concentrate well.*)
- *"Vielleicht haben Sie auch zu tief gebückt"* (baseline) (*Maybe you have bent over too far.*)

In German, both verbs (*"concentrate"* and *"bend over"*) require the reflexive pronoun *"sich"* (themselves or yourself), which, however, the tool omitted.

- Incorrect statements and advice:

- *"Nehmen Sie Ihren Helm ab"* (model 2) (*take off your helmet*). In this text about first aid, the tool erroneously capitalized the pronoun *"Ihren"*, which therefore refers to the second person singular instead of the third person singular. The correct spelling would be lowercase (*"ihren"*).
- *"Sie können die Pille auch in der Schwangerschaft nehmen"* (model 1) (*You can also take birth control pills during pregnancy.*)
- *"Und Sie sollten alles tun, was Ihren Gelenken schadet."* (model 1) (*And you should do anything that harms your joints*). The verb of the source text *"meiden"* (avoid) was translated with its antonym *"tun"* (do). Thus, the reader is even given harmful advice.
- *"Bei etwa 14 Prozent der Patienten [...] ist die Herz-Kranz-Gefäße verengt."* (model 1) (*In about 14 percent of patients, the coronary arteries is narrowed*). In this case, not only the verb form is incorrect (singular instead of plural), but the tool also failed to translate the negated statement of the source text (*"findet sich [...] keine Verengung der"*

Koronargefäße)” (*no narrowing of the arteries is found*).

- Incorrect explanations of technical terms: *”Die Zeit, in der man krank ist, nennt man Inkubationszeit.“* (baseline) (*The period during which one is sick is called incubation period*). This is incorrect because the incubation period is the time between the infection and the manifestation of symptoms.
- Wrong relation: Source text: *”Deshalb ist hier unbedingt ein Arztbesuch angeraten. Auch wenn die Symptome länger als drei Tage anhalten [...] wird der Besuch beim Arzt unumgänglich”* (*Therefore, seeing a doctor is strongly recommended here. Also, if the symptoms persist for more than three days, a visit to the doctor is inevitable.*) vs. Target text: *”Bei diesen Menschen kann eine Lebens-Mittel-Vergiftung schwerer verlaufen. Deshalb sollten Sie bei diesen Anzeichen sofort zum Arzt gehen; Die Beschwerden dauern länger als 3 Tage.”* (model 1) (*In these individuals, food poisoning can progress more severely. Therefore, if you experience these symptoms, you should see a doctor immediately: The symptoms last longer than 3 days.*). In the source text, the word *”hier”* (here) refers to vulnerable groups of people; however, it is erroneously translated with *”symptoms”*. As a result, the target text states that these individuals should only see a doctor when they experience one of the following symptoms, while the source text indicates that vulnerable people have to see a doctor in any case.
- Homophonic but not homographic words are not correctly selected: *”Dann 7 Sie den Saft durch ein Tuch oder einen Kaffeefilter.”* (model 2) (*Then strain the juice through a cloth or a coffee filter*). In German, the word *”sieben”* is both a verb (strain) and a number (7).

Correctness is not yet present for the different systems under study to the extent that texts would be usable without post-editing. The human translation corpus does not have such errors, but has a high degree of correctness.

3.3.2 Readability

Comparing the comprehensibility of the different corpora revealed that, as expected, the source texts were the least comprehensible texts (mean: 10.46, SD: 2.76). Model 2 had the highest comprehensibility, with a mean HIX value of 19.5 (SD: 0.76). While this is a slight improvement compared to the baseline model (mean: 19.15, SD: 0.49), Figure 1 shows that the HIX value for model 1 (mean: 17.71, SD: 1.41) was considerably lower than that of the baseline model, i.e., based on the HIX, the model’s comprehensibility was not improved. However, as seen from the boxplot, human translations also yielded a lower HIX value (mean: 17.74, SD: 1.67) than the baseline model, and both the human and the model 1 translations reveal a much greater variation in the HIX values than the baseline and model 2 translations. While from the model 1 translations, only 93% of the texts, and from the human translations, only 83% of the texts reached the predefined Plain German benchmark, all of the baseline and model 2 texts could be classified as Plain German texts. However, when interpreting the HIX value, it should be kept in mind that this is only a quantitative analysis that focuses only on comprehensibility features on the text surface (i.e. overt complexity) and the textual level is mainly ignored. For this reason, HIX values only represent a starting point for the analysis and it has to be complemented by a qualitative analysis (e.g. Section 3.3.1).

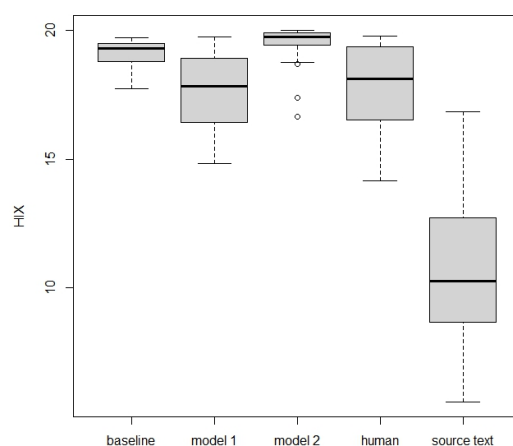


Figure 1: HIX values of the three machine translations, the human translations, and the source texts.

3.3.3 Syntactic Complexity

As seen from Figure 2, human translations contain the least number of complex syntactic con-

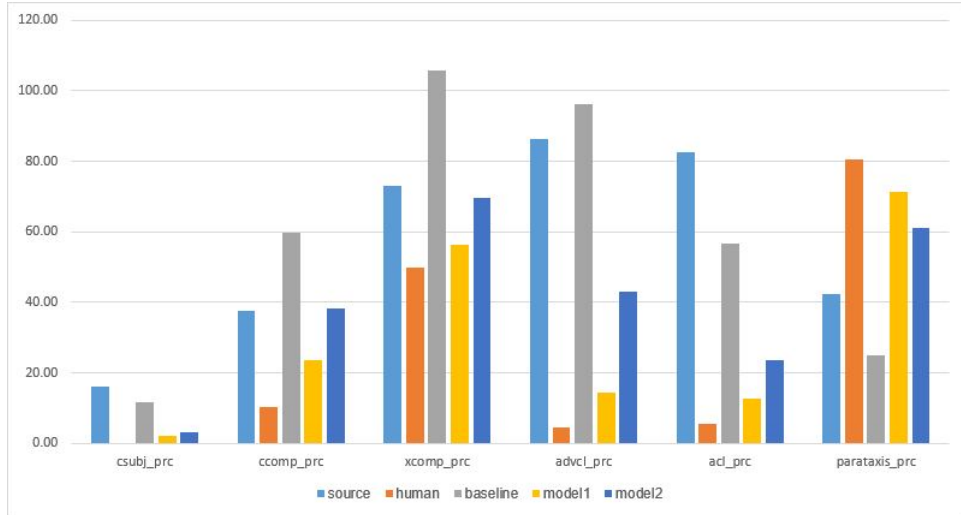


Figure 2: Distribution of syntactically complex dependency relations in the source texts, human and machine translations (normalised frequencies per 10000).

structions, except for parataxis.

In machine translation outputs, we observe the following pattern: model 1 contains the least number of complex syntactic constructions, followed by model 2. Here again, the only exception is the distribution of the parataxis constructions. Translations with the baseline model are much more complex in terms of syntax if compared to the other two systems. Remarkably, for some structures, they are even more complex than the source texts.

3.3.4 Automatic evaluation measures

In the final stage, we examined the text simplification metric SARI. Figure 3 displays box-

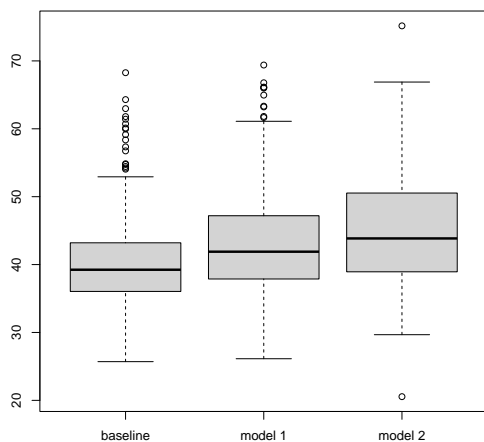


Figure 3: SARI score of the machine translation output from the baseline model, model 1 and model 2.

plots, comparing the SARI of the machine translation output from the baseline model, model 1 and

model 2. Higher SARI values indicate better machine translated outputs.

The system utilized in our analysis achieves an average SARI score of 40.61 (SD: 6.78) for the baseline model, 43.49 (SD: 7.84) for model 1 and 45.13 (SD: 8.15) for model 2. All models are therefore in line with with state-of-the-art text simplification models reported by Sheang and Saggion (2021).

3.3.5 Translation Comparison

We now summarise the results of all evaluation criteria for the system outputs. Table 1 illustrate the system ranking depending on the used criteria.

system	corr	HIX	synt	SARI
baseline	3	2	3	3
model 1	2	3	1	2
model 2	1	1	2	1

Table 1: System ranking according to the evaluation criteria: corr=correctness, HIX=readability, synt=syntax, and SARI=text simplification.

As seen from the table, the worst outputs according to our evaluation criteria are found with the baseline system. The results for models 1 and 2 vary depending on the evaluation criteria. For instance, model 1 performs better in terms of syntax, as its outputs reveal not so many complex syntactic constructions. This system seems to be very close to human translations as well. However, many of the texts translated with model 1 are not correct. In terms of correctness, as well as readability scores and text simplification scores, model 2 is the win-

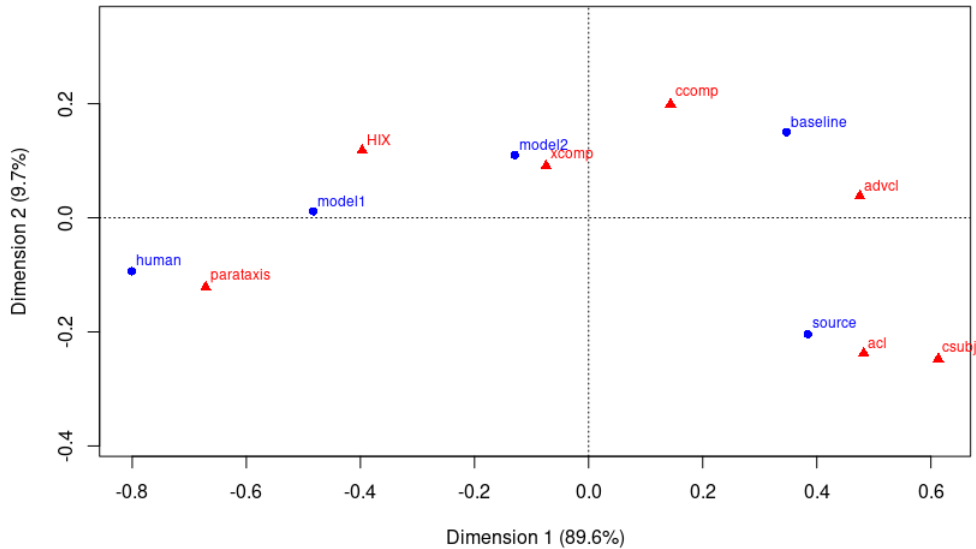


Figure 4: Correspondence analysis based on HIX scores and distribution of syntactic structures.

ner.

For the comparison of the machine-translated outputs with the human translations and sources, we only use HIX scores and distributions of the complex syntactic structures. The results of the correctness analysis are not numeric. The resulting two-dimensional graph is shown in Figure 4.

The most obvious information we can obtain from this graph is that the difference is most strongly pronounced along the x-axis between the two groups of subcorpora: source texts and translations with the baseline model on the right side vs. human translations and outputs of the two other models (model 1 and model 2) on the left side. This difference is considerable, as the dimension along the x-axis explains a very high proportion (89.6%) of the data variance. We also see that translations with model 1 are the closest to the human translations. On the y-axis, we see a separation between human- vs. machine-authored texts. However, this difference is not big, as it explains only 9.7% of the data variance in our dataset.

4 Discussion and Future Work

The present paper evaluates three different models of a machine translation system for translating medical texts into Plain German. It covers one of the first steps towards the implementation of these tools in accessible health communication in Germany and it discusses first methodological ap-

proaches, which we intend to expand on in further research.

Model 2 seems to achieve the best results according to most criteria. At the same time, model 1 seems to be more similar with the human translations at hand. While in terms of syntactic complexity and text readability, the models yielded promising results, the evaluation of the correctness revealed severe misinformation for all three models, the consequence being that the texts cannot be safely used by end users. At the same time, the tool under analysis can be used as a CAT tool for professional translators and content providers with an expertise in Plain Language and post-editing. As our study has clearly revealed that so far machine translated Plain Language texts cannot do without post-editing, but need intensive revision, professional post-editing competences are more important than ever. This means that translators and experts working on machine translated text must be trained to detect and correct different types of errors, especially those that are critical for user safety. In further steps, a guide for post-editing in intralingual translation will be developed, exposing the necessary competences and factors to be considered when using machine translation into Plain German.

In a next step, the machine translation system will now be integrated into the editorial workflow of the *Apotheken Umschau* on a trial basis. This

practice test will serve to assess the time and effort that is needed to post-edit the machine translated texts. Adding machine translation to the editorial process could optimize the process and addresses the gap between the need for texts in Plain German and the lack of professional translators (see Section 1). The metrics from the practice test will be particularly interesting because the tool will only be permanently integrated into the workflow if the time and effort for post-editing the output is lower than for translating from scratch. Therefore, these metrics will determine the final decision for or against adapting the status quo.

In our future research, we will conduct a thorough analysis and classification of the various error types found in the machine translated texts. For example, we plan to investigate specific linguistic phenomena, such as the translation of compound words.

In addition, we also want to test and compare the output from both SUMM AI and other state-of-the-art systems to investigate which of the currently available systems is most suitable for intralingual translation into Plain German, in general and for specific subjects and text types. These systems include both freemium tools that offer both free and paid plans, such as ChatGPT and Google Gemini, and commercial tools, such as Klartext St. Pauli, capito digital and T2K (text2knowledge). In these future studies, we also plan to use other text types and texts from other domains, so that we are able to compare not only different tools but also different datasets.

References

- Ahrens, Sarah, Rebecca Schulz, Janina Kröger, Sergio Hernández Garrido, Loraine Keller, and Isabel Rink. 2022. Accessible communication and health literacy. *Accessibility–Health Literacy–Health Information: Interdisciplinary Approaches to an Emerging Field of Communication*, 13:9.
- Anschütz, Miriam, Joshua Oehms, Thomas Wimmer, Bartłomiej Jezierski, and Georg Groh. 2023. Language Models for German Text Simplification: Overcoming Parallel Data Scarcity through Style-specific Pre-training. *arXiv preprint arXiv:2305.12908*.
- Blechschmidt, Anja. 2021. Health literacy and multimodal adapted communication. *New Approaches to Health Literacy: Linking Different Perspectives*, pages 65–82.
- Bredel, Ursula and Christiane Maaß. 2016. *Leichte Sprache: Theoretische Grundlagen, Orientierung für die Praxis*. Dudenverlag.
- De Marneffe, Marie-Catherine, Christopher D Manning, Joakim Nivre, and Daniel Zeman. 2021. Universal dependencies. *Computational linguistics*, 47(2):255–308.
- Deilen, Silvana, Sergio Hernández Garrido, Ekaterina Lapshinova-Koltunski, and Christiane Maaß. 2023a. Using ChatGPT as a CAT tool in Easy Language translation. In *Proceedings of the Second Workshop on Text Simplification, Accessibility and Readability (TSAR)*, RANLP, Varna, Bulgaria. ACL.
- Deilen, Silvana, Silvia Hansen-Schirra, Sergio Hernández Garrido, Christiane Maaß, and Anke Tardel. 2023b. Emerging fields in easy language and accessible communication research. In *Emerging Fields in Easy Language and Accessible Communication Research*, pages 9–15. Springer.
- Deilen, Silvana. 2021. *Optische Gliederung von Komposita in Leichter Sprache. Blickbewegungsstudien zum Einfluss visueller, morphologischer und semantischer Faktoren auf die Verarbeitung deutscher Substantivkomposita*. Frank & Timme.
- Greenacre, Michael J. 2007. *Correspondence Analysis in Practice*. Chapman & Hall/CRC, Boca Raton.
- Gutermuth, Silke. 2020. *Leichte Sprache für alle?: eine zielgruppenorientierte Rezeptionsstudie zu Leichter und Einfacher Sprache*, volume 5. Frank & Timme GmbH.
- Hansen-Schirra, Silvia, Walter Bisang, Arne Nagels, Silke Gutermuth, Julia Fuchs, Liv Borghardt, Silvan Deilen, Anne-Kathrin Gros, Laura Schiffl, and Johanna Sommer. 2020a. Intralingual translation into easy language—or how to reduce cognitive processing costs. *Easy Language Research: Text and User Perspectives*. Berlin: Frank & Timme, pages 197–225.
- Hansen-Schirra, Silvia, Jean Nitzke, Silke Gutermuth, Christiane Maaß, and Isabel Rink. 2020b. Technologies for translation of specialised texts into easy language. *Easy Language Research: Text and User Perspectives*. Berlin: Frank & Timme, pages 99–127.
- Kopp, Tobias, Amelie Rempel, Andres P. Schmidt, and Miriam Spieß. 2023. Towards machine translation into easy language in public administrations: Algorithmic alignment suggestions for building a translation memory. In Deilen, Silvana, Silvia Hansen-Schirra, Sergio Hernández Garrido, Christiane Maaß, and Anke Tardel, editors, *Emerging Fields in Easy Language and Accessible Communication Research*, pages 371–406. Frank & Timme, Berlin.
- Maaß, Christiane and Isabel Rink. 2021. Translating legal texts into Easy Language. *J. Open Access L.*, 9:1.

- Maaß, Christiane and Laura Marie Schwengber. 2022. Easy Language and Plain Language in Germany. *Rivista internazionale di tecnica della traduzione= International Journal of Translation*.
- Maaß, Christiane, Isabel Rink, Silvia Hansen-Schirra, Camilla Lindholm, and Ulla Vanhatalo. 2021. Easy language in Germany. *Handbook of Easy Languages in Europe*, 8:191.
- Maaß, Christiane. 2020. *Easy Language–Plain Language–Easy Language Plus: Balancing Comprehensibility and Acceptability*. Frank & Timme.
- Maddela, Mounica, Fernando Alva-Manchego, and Wei Xu. 2021. Controllable text simplification with explicit paraphrasing. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3536–3553, Online, June. Association for Computational Linguistics.
- Martin, Louis, Éric de la Clergerie, Benoît Sagot, and Antoine Bordes. 2020. Controllable sentence simplification. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4689–4698, Marseille, France, May. European Language Resources Association.
- Ondov, Brian, Kush Attal, and Dina Demner-Fushman. 2022. A survey of automated methods for biomedical text simplification. *Journal of the American Medical Informatics Association*, 29(11):1976–1988.
- R Core Team, 2017. *R: A Language and Environment for Statistical Computing*. Vienna, Austria.
- Rink, Isabel. 2019. *Rechtskommunikation und Barrierefreiheit: Zur Übersetzung juristischer Informations- und Interaktionstexte in Leichte Sprache*. Frank & Timme.
- Saggion, Horacio. 2017. Applications of automatic text simplification. In *Automatic Text Simplification*, pages 71–77. Springer.
- Säuberli, Andreas, Sarah Ebling, and Martin Volk. 2020. Benchmarking data-driven automatic text simplification for German. In *Proceedings of the 1st workshop on tools and resources to empower people with reading difficulties (READI)*, pages 41–48.
- Scarton, Carolina and Lucia Specia. 2018. Learning simplifications for specific target audiences. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 712–718, Melbourne, Australia, July. Association for Computational Linguistics.
- Schaeffer, Doris, Dominique Vogt, Eva-Maria Berens, and Klaus Hurrelmann. 2017. *Gesundheitskompetenz der Bevölkerung in Deutschland: Ergebnisbericht*. Bielefeld: Universität Bielefeld, Fakultät für Gesundheitswissenschaften.
- Schaeffer, Doris, Klaus Hurrelmann, Ullrich Bauer, and Kai Kolpatzik. 2018. Nationaler Aktionsplan Gesundheitskompetenz. *Die Gesundheitskompetenz in Deutschland stärken*. Berlin: KomPart, 10:0418–1866.
- Schaeffer, Doris, Eva-Maria Berens, Svea Gille, Lennert Griese, Julia Klinger, Steffen de Sombre, Dominique Vogt, and Klaus Hurrelmann. 2021. Gesundheitskompetenz der Bevölkerung in Deutschland vor und während der Corona Pandemie: Ergebnisse des HLS-GER 2. Technical report, Universität Bielefeld, Interdisziplinäres Zentrum für Gesundheitskompetenzforschung.
- Sheang, Kim Cheng and Horacio Saggion. 2021. Controllable sentence simplification with a unified text-to-text transfer transformer. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 341–352, Aberdeen, Scotland, UK, August. Association for Computational Linguistics.
- Sørensen, Kristine, Stephan Van den Broucke, James Fullam, Gerardine Doyle, Jürgen Pelikan, Zofia Slonska, Helmut Brand, and (HLS-EU) Consortium Health Literacy Project European. 2012. Health literacy and public health: a systematic review and integration of definitions and models. *BMC public health*, 12:1–13.
- Spring, Nicolas, Marek Kostrzewa, David Fröhlich, Annette Rios, Dominik Pfützte, Alessia Battisti, and Sarah Ebling. 2023. Analyzing sentence alignment for automatic simplification of German texts. In *Emerging Fields in Easy Language and Accessible Communication Research*, pages 339–369. Springer.
- Xu, Wei, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. Optimizing statistical machine translation for text simplification. *Transactions of the Association for Computational Linguistics*, 4:401–415.

Using Machine Learning to Validate a Novel Taxonomy of Phenomenal Translation States

Michael Carl
Kent State University, USA
mcarl16@kent.edu

Sheng Lu
TU Darmstadt, Germany
boblu@ruc.edu.cn

Ali Al-Ramadan
Kent State University, USA
aalramd@kent.edu

Abstract

We report an experiment in which we use machine learning to validate the empirical objectivity of a novel annotation taxonomy for behavioral translation data. The HOF taxonomy defines three translation states according to which a human translator can be in a state of Orientation (O), Hesitation (H) or in a Flow state (F). We aim at validating the taxonomy based on a manually annotated data-set that consists of six English-Spanish translation sessions (approx 900 words) and 1813 HOF-annotated Activity Units (AUs). Two annotators annotated the data and obtain high average inter-annotator accuracy 0.76 (kappa 0.88). We train two classifiers, a Multi-layer Perceptron (MLP) and a Random Forest (RF) on the annotated data and tested on held-out data. The classifiers perform well on the annotated data and thus confirm the epistemological objectivity of the annotation taxonomy. Interestingly, inter-classifier accuracy scores are higher than between the two human annotators.

1 Introduction

Translation is considered to involve complex and non-linear cognitive processes (Krings, 2001). Understanding the intricacies of the temporal dynamics of these processes is a fundamental aspect in Translation Process Research (TPR).

Various approaches have been proposed over the past 40 years to understand the distinct phases and

mental states experienced by translators (Jakobsen, 2017). Starting with Think-Aloud Protocols in the 1980s, in which translators comment their own translation behavior during their translations (Königs, 1987; Krings, 2001), the field of enquiry has moved towards less invasive technologies, that is, keystroke logging and eye tracking (Hvelplund, 2016; Carl et al., 2016). The recordings of these logging tools make it possible to assess the flow of translation in a seamless way and to investigate how translations evolve in time; where translators type smoothly, where they get stuck, and where they search for (external) resources, etc.

One approach to analysing the translation process has been to segment the behavioral Translation Process Data (TPD) into processing units (Alves and Vale, 2009; Schaeffer et al., 2016). But how these segments should be defined and what they exactly represent has been a topic of continuous exploration and debate. The assessment of the translation rhythm (aka "Pause Analysis" (Kumpulainen, 2015; Muñoz and Apfelthaler, 2022)) has provided valuable insights into translation patterns as produced by more or less experienced translators (Jakobsen, 2011), for different levels of text complexity (Hvelplund, 2016), for different translation goals (Zou et al., 2022b), post-editing behavior (Jia et al., 2019) and also for spoken translation (e.g., interpretation, sight translation, (Zou et al., 2022a)). The underlying assumption has been that longer keystroke pauses are indicative of more challenging translations, while stretches of smooth typing can be observed when there are no/less translation hurdles or difficulties (Lacruz et al., 2014). However, determining an exact pause threshold to differentiate these phenomena remains a challenge. Many studies (Krings, 2001; O'Brien, 2006; Kumpulainen, 2015; Vieira, 2016,

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

among others) present varying segmentation methods with different pause thresholds, ranging from 300ms to five seconds or more. These strands of research rely on deterministic fragmentation techniques to segment the key logging and eye tracking data into Translation Units (TUs, (Alves and Vale, 2009; Carl and Schaeffer, 2017)) or Activity Units (AUs, (Hvelplund, 2016; Schaeffer et al., 2016)). However, these approaches lack intuitive labeling and thus make it difficult to derive a comprehensive understanding of the complex nature of the translation process and how it unfolds over time. Some researchers suggest a hierarchical process model (Schaeffer and Carl, 2013) and others (Muñoz and Apfelthaler, 2022; Dragsted, 2010) advocate a translator-specific fragmentation of the TPD into processing units depending on the translators' typing speed.

Combined, this suggests that human translation processes are embedded in a hierarchical mental architecture, encapsulating various processing strata. A hierarchical approach to understanding translation has the potential to offer a nuanced understanding of translators' behavior and strategies on various interacting levels of analysis. In order to advance this project, a novel higher-level segmentation taxonomy was introduced in (Carl et al., 2024) that fragments the TPD in three broad phenomenal translation states, Hesitation (H), Orientation (O) and Translation Flow (F). The HOF taxonomy assumes that behavioral traces of these three states can be observed in the TPD and that translators can be at any one point in time in only one of the three states.

In previous work (Carl et al., 2024) we have annotated a small corpus with HOF translation states. The corpus is publicly available as part of the CRITT TPR-DB¹. The HOF annotation corpus provides a layer of manual annotation, introducing segment labels of an assumed phenomenal layer of translation processes, suited to analyse the hierarchical embedding of translation processes.

However, the HOF taxonomy, capturing qualities of conscious translator experience, is entirely new territory and the generalizability and validity of the annotation schema is still unclear. In this paper, we therefore conduct further investigation to assess whether and to what extent this new taxonomy is epistemologically valid — that is, we want

to investigate to what extent different annotators might agree the HOF states to represent a phenomenal translation "reality". While a manual annotation has shown a varying amount of agreement between two annotators (Kappa 0.37 and 0.88, (Carl et al., 2024)), in this paper we use ML techniques to further validate the consistency of the taxonomy.

AI and ML techniques can be used in various ways and for different purposes. Mollo (2024), for instance, enumerates a few scenarios where AI can be used as: AI-as-engineering (industrial and commercial projects, as e.g., MT systems), AI-as-psychology to improve our understanding of biological intelligence, AI-as-idea or AI-as-recreation for recreating biological intelligence in artificial systems. AI can also be used for "exploring intelligence spaces" so as to uncover new forms of intelligence that are different from human intelligence or to uncover algorithms (Zhong et al., 2023).

In this paper we use ML techniques to investigate the "epistemological objectivity" (Searle, 1998; Searle, 2017, see also section 5 for a discussion) of the HOF annotation schema. That is, we are interested in verifying whether ML can reproduce the results of our manual annotations to a similar amount of accuracy. We assume that if the trained models performs well on the classification task, it confirms the objectivity of the annotation taxonomy used to create the training data. Conversely, poor model performance would indicate issues with the annotation taxonomy.

For instance, it might be the case that, even though two annotators agree in their annotation label, they might be biased by some intuition, cultural or otherwise un-observable features which may not be accessible to the ML technology. However, if ML reaches similar results of accuracy as the inter-annotator agreement indicates, we take the annotation taxonomy to implement reproducible and epistemological objective annotation criteria. That is, as ML lacks subjective, personal or cultural influences in the process annotation process (i.e., HOF state labeling), high accuracy on held-out testing data may be an indicator of stable results with minimum bias,

In section 2 we describe the data and the manual annotation process of the reference (training) data. Section 3 describes our implementation of two classifiers — Multi-Layer perception (MLP) and Random Forest (RF) — while section 4 reports our training and evaluation on a set of 1813

¹The annotations can be downloaded from here http://critt.as.kent.edu:3838/public/State_Annotation_Phases.zip

HOF-annotated AUs. We report higher precision and recall between the classifier and annotator as compared to inter-annotator agreement between the two annotators. Section 5 concludes with a discussion on Searle’s notion of “ontological subjectivity” and “epistemological objectivity” and their relation to the evaluation of our HOF states.

2 Activity Units, Translation Units and Translation States

The translation process can be conceptualized as a dynamic flow of mental processes marked by information input as gathered through eye movements and textual output in the form of keystrokes or mouse movements. As translators navigate the source text (ST) and produce the target text (TT), their behavior is influenced by numerous factors. To better understand these processes, different approaches have been pursued that fragment translation-behavioral data (keystrokes and gaze data) into processing units. Figure 1 shows a progression graph that depicts three ways of segmenting the approximately 28 seconds of the plotted translation session. TUs (Carl and Kay, 2011; Alves and Vale, 2009), indicated in the top in Figure 1, are characterized by a typing pause (a blank space in top line of the Figure) followed by a typing burst (or Production Unit, PU) indicated as grey boxes. AUs (Schaeffer et al., 2016; Hvelplund, 2016) are constructed based on the coordination of gaze activities and typing behavior and are marked at the bottom in Figure 1 in different colors. Three distinct HOF translation states are indicated with black dotted lines (Carl et al., 2024). Boundaries of translation states coincide with AU boundaries, so that sequences of AUs can be used to fragment the TPD into HOF translation states. This section explores these constructs and the three annotations, emphasizing their significance and interplay.

2.1 Translation Units

TUs segment the continuous stream of translation activities (keystrokes) into stretches of typing and pausing. They capture the translator’s perception and actions, indicating the challenges they encounter during the translation process (Malmkjær, 1998). In Figure 1, TUs appear as successive pauses and typing bursts of fluent production (or PUs). The pauses that occur between PUs are taken to be indicators of elevated translation effort,

as it is assumed that during those breaks translators engage in reflection or (mental) search (Dragsted, 2010). Sequences of TUs have been used to compute pause-word-ratio (Lacruz et al., 2014) as indicators of cognitive effort. However, TUs often lack the granularity to explain precisely what occurs during the pauses. Moreover, they do not differentiate whether a translator is directing their focus towards the ST or the TT during these intervals (Schaeffer et al., 2016).

AU	AU activity	Color	Effort	Effect
T1	ST reading	Blue	+	-
T2	TT reading	Green	+	-
T4	translation production	Yellow	-	+
T5	ST reading with concurrent production	Red	-	+
T6	TT reading with concurrent production	Dark Green	-	+
T8	no observed behavior for more than one second	Black	+	-

Table 1: Types of AUs, color code in Figure 1 and levels of translational effect and cognitive effort.

2.2 Activity Units

AUs provide a more fine-grained view on the translation process, focusing on the coordination of the translator’s eyes and hands. It addresses some of the inherent limitations of TUs. In our data, we categorize AUs into six types as presented in Table 1 (Carl et al., 2016). They provide more detailed insights into how translators engage in various aspects of the translation process. The classification of AUs is based on whether translators are actively involved in translation production, reading the ST or TT, or simultaneously reading and writing. As shown in Table 1, each type of AU can be associated with a degree of translational effects (typing activities) and cognitive effort (i.e., gazing). For instance, an AU of type T1 indicates ST reading which results in low levels of effects (no translation is typed) but higher amounts of cognitive effort (mental resources are allocated). In simpler terms, it means that the translator primarily focuses on understanding the ST, with minimal to no simultaneous translation work, as depicted in Figure 1.

2.3 HOF Translation States

HOF Translation States offer insights into qualities of the translator’s experience. Carl et al (2024) distinguish between three translation states: A state

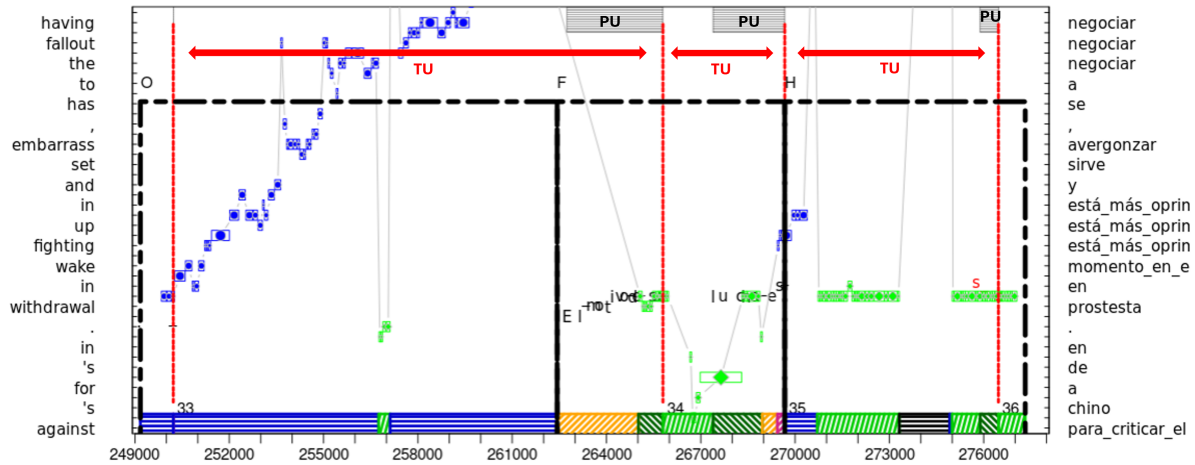


Figure 1: Progression graph of a small snippet of the translation session (BML12/P02.T3). Production time in milliseconds is indicated in the horizontal axis. Vertical axis refers to the ST on the left side and the TT on the right side. The blue dots and green diamonds represent eye movements on the ST and TT respectively. The black and red characters are insertion and deletion respectively. AUs are marked as colored bars on the bottom, TUs are indicated with red lines, and PUs as gray boxes in the top of the graph. Translation States are sequences of AUs, indicated by black dashed boxes, labeled O, F, H. The graph represents a segment of approximately 28 seconds (249.000ms - 277.000ms) of an English-to-Spanish translation.

of orientation (O) refers to the translator’s behavior when feeling the need to get acquainted with the source text (ST). It is characterized by linear, forward-reading behavior of the ST. The Flow state (F) represents a phase in which the translator is immersed in translation production, generating the TT with ease and minimal interruption. It is marked by fluent translation production with minimal reading ahead and short pauses. A state of Hesitation (H) emerges out of surprise, where unexpected challenges prompt the translator to revise and re-read. This state indicates moments of uncertainty or cognitive challenge, signifying areas where the translator is challenged with complexities in the source text or struggles to find suitable translations. These distinct translation states are annotated in the progression graph in Figure 1, exemplifying associated typical behavioral correlates.

2.4 Empirical Data

We use a set of six translation sessions from the CRITT TPR-DB that were previously annotated with HOF translation state labels (Carl et al., 2024). The CRITT TPR-DB (Carl et al., 2016) is a collection of currently more than 5000 translation sessions, amounting to hundreds of hours of TPD, that is compiled into a consistent publicly available database. The CRITT TPR-DB is extensively documented in numerous publications and summary tables with more than 300 product and process fea-

tures are available in a compiled form.²

In this study we use six English-to-Spanish translation sessions from BML12³. The BML12 study consists of 184 translation sessions by Spanish translation students that were recorded in 2012 in Copenhagen and in Spain (Barcelona). The HOF annotation taxonomy was developed based (among others) on six BML12 sessions and annotated in 2022 by two advanced (Chinese and Japanese) translators. A special purpose interface was used to annotate the translation sessions, similar to Figure 1. The annotation process is described in detail (Carl et al., 2024). The six annotated sessions consist in total of 42 segments (sentences) with 854 source words. The translations of these 42 segments resulted in 1813 AUs which were annotated with HOF labels. In this study we used the 1813 HOF-annotated AUs for training and evaluating two classifiers.

²See the CRITT website <https://sites.google.com/site/centrereinnovation/tpr-db>. The TPD can be downloaded free of charge from sourceforge <https://sourceforge.net/projects/tprdb/>, an introduction to the usage and a free trial account is provided here <https://sites.google.com/site/centrereinnovation/tpr-db/getting-started>

³The MultiLing data and BML12 study is described: https://sites.google.com/site/centrereinnovation/tpr-db/public-studies#h.p_iVVuCQOHJx20

2.5 Manual Annotation

As reported in (Carl et al., 2024), the manual annotation involved five phases: Phases 1, 2, and 4 were trial annotations and are not considered here. In Phase 3, 1288 AU were annotated, but the absence of a structured approach resulted in a Kappa score of 0.37, indicating a moderate agreement between the two annotators (see Table 2). In Phase 5, a structured approach with a decision tree and guidelines was defined (Carl et al., 2024), resulting in a significant improvement in inter-annotator agreement, as shown in the high Kappa score of 0.88. Table 2 shows the inter-rater accuracy and Kappa scores along with the number of AUs used in Phase 3 and 5.

Phase	Total AUs	Kappa	Accuracy
3	1288	.37	.66
5	525	.88	.93

Table 2: Kappa scores for Phases 3 and 5 of annotation. Phase 3 involves five sessions, while Phase 5 involves session P04_T2 of the BML12 study. The average accuracy for all 1813 annotations is .74.

Furthermore, Table 3 offers a breakdown of the number of AUs for each of the three translation states in Phases 3 and 5, and for both annotators (Y and T). There is a noticeable shift in the distribution of annotated AUs across these states between the two phases. In Phase 3, the difference in the numbers of AUs annotated by T and Y across the three states suggests distinct interpretations of the states. This is the reason for the low Kappa score of 0.37 and low Accuracy of 0.66 in Table 2. The elaboration of a decision tree and annotation guidelines prior to Phase 5 clearly leads to a better alignment between the two annotators, evidenced not only by their closely matching counts across states but also by the high Kappa score of 0.88 (Accuracy 0.93). In our experiments, we use annotations from Phases 3 and 5 as a training/test corpus in section 4. Given the amount of coordination and mutual adjustment of the two annotators, we decided to corroborate the empirical objectivity of the annotation schema using MT.

3 Training Translation State Classifiers

In this section we describe two classifiers that were used to assess the annotated translation states. While the variation of inter-annotator agreement, as reported in Table 2, indicates that annotators

State	Phase 3		Phase 5	
	Y	T	Y	T
H	403	331	216	217
O	275	108	51	55
F	610	849	258	253
total	1288	1288	525	525

Table 3: Number of AUs in the two AU annotation phases for the two annotators T and Y.

are able to learn and agree on annotation guidelines and to generalize and reproduce the underlying concepts, this does not necessarily mean that those generalizations can also be learned and reproduced by a ML classifier. If, however, ML techniques can reproduce manual annotations with high accuracy we can be more certain about the "epistemological objectivity" (Searle, 2017) of the annotation schema. Besides, once a classifier is trained, we will also be able to automatically annotate new data. Therefore, in this study we only describe an evaluation of the trained classifier on the manually annotated data and leave a full-blown analysis on a large corpus for future research.

3.1 Multi-layer Perceptron

An MLP is a supervised simple feed-forward neural network. It consists of multiple layers, and each layer is fully connected to the following one. Figure 2 shows the structure of a two-layer MLP.

For our task, there are 34 cells in the input layer, and each of them corresponds to one of the 34 AU features (see Appendix A). There are 3 cells in the output layer corresponding to 3 state labels {H, O, F}. The output is a set of probabilities, each of which represents the probability that an input is classified as a certain state. The final prediction is the state with the highest probability.

We implemented MLP classifiers using `sklearn`⁴. We set the following parameters in `MLPClassifier`:

- `hidden_layer_sizes`: number of neurons in hidden layers
- `batch_size`: size of mini-batches
- `max_iter`: maximum number of iterations
- `learning_rate_init`: learning rate used
- `random_state`: random seed
- `solver`: weight optimizer

⁴See <https://scikit-learn.org>.

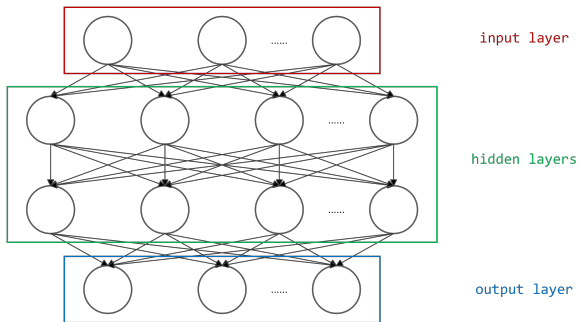


Figure 2: The structure of a two-layer MLP.

3.2 Random Forest

We implemented an RF classifier using `sklearn` (see footnote 4). RF is a supervised ML method that trains a model on the annotated (‘gold’) data. An RF is a set of decision trees where the output of the classifier is the class selected by most trees (majority vote). In our trials we set this number to 500. RFs are said to be rather robust with respect to variations in the data, as the entry point to each of the decision trees varies probabilistically, while the algorithm averages over the differences. New, unseen data can, therefore, be classified reliably. Another advantage of RFs is that the importance of the features can be ranked, which may provide helpful insights into feature design. In Appendix B we provide the ranking of feature importance for the 34 AU features and the two models trained on the annotated Y and T data.

4 Evaluation of Classifiers

The training of AU-to-state classifiers is based on six annotated sessions from Phase 3 and 5 with 1813 datapoints (AUs), as shown in Table 3.

We trained classifiers based on two backbone models: MLP and an RF. The annotated data are split into 70% for training (1269 data points), and 30% for testing (544 data points). We used the 34 features that are described in the Appendix A in Table 9 for the classifiers.

In a 10-fold cross validation with RF and MLP we get best average accuracy values of 0.85 for the two classifiers. Table 4 shows the best and average performance of classifiers based on the two models. We observe the best accuracy and F1-score for the RF classifiers and annotator T.

4.1 Multi-layer Perceptron

For the MLP classifier, we used the Adam optimizer (Kingma and Ba, 2014) with an alpha value

		Best		Average	
		acc.	F1	acc.	F1
T	MLP	85	69	67	52
	RF	85	75	78	64
Y	MLP	79	73	66	51
	RF	83	76	75	57

Table 4: Best and average accuracy (acc.) and F1-score (F1, in percentage) for the RF and MLP classifiers for both annotators.

of $1e-5$. The MLP architecture consisted of three hidden layers with 400, 200, and 400 units in size, respectively. We set the random state to 1 for reproducibility. The classifiers were trained separately on the data annotated by annotators Y and T. We also used the `StandardScaler` from `sklearn`, which standardizes features by removing the mean and scaling to unit variance.

4.2 Scaling Data

It is worth noting a difference in the performance of the MLP classifiers when trained on scaled data vs. non-scaled data. The effect of `StandardScaler` to the performance of the MLP classifier is significant. When the MLP model is trained on unscaled data, its average accuracy drops significantly. However, the scaler does not impact the results of the RF classifier.

MLP	State	Prec.	Rec.	F1	Support
T	H	0.77	0.79	0.78	164
	O	0.72	0.69	0.70	45
	F	0.91	0.90	0.91	335
Y	H	0.76	0.74	0.75	204
	O	0.66	0.59	0.63	96
	F	0.81	0.86	0.83	244

RF	State	Prec.	Rec.	F1	Support
T	H	0.82	0.83	0.82	168
	O	0.72	0.48	0.57	48
	F	0.90	0.94	0.92	328
Y	H	0.84	0.84	0.84	189
	O	0.84	0.73	0.78	93
	F	0.88	0.92	0.90	262

Table 5: Precision (Prec.), Recall (Rec.), and F1-score (F1) for MLP (top) and RF (bottom) classifiers trained on 1269 AU annotations and evaluated on a test set of 544 AU annotations for both annotators.

4.3 Precision and Recall

A fine-grained assessment of the classification report in Table 5 reveals that values for Precision, Recall and F1-score are differently distributed for the three States. State F has highest precision and recall values for both annotators Y and T and for both classifiers. As noted in (Carl et al., 2024), it is comparatively easy to detect this state in the behavioral data and it also has the best inter-rater agreement, as discussed in section 2 and Table 3. States O and H seem to be more difficult to separate and may be easily confused. Note that that Table 5 shows this to be the case for both annotators, T and Y.

4.4 Comparing Y and T labels

Table 6 shows two confusion matrices between the Y annotations and the RF predictions (on the left) and the T annotations (on the right). As the two matrices show, as well as indicated in Tables 3 and 5, the distribution of states are unequally distributed. There are almost three times more AUs with F label, as compared to O states.

A large number of states H seem to be classified as F, which indicates that more distinctive features might need to be developed, so as to better distinguish between states F and H.

True Y-labels	RF Predictions			True T-labels		
	H	O	F	H	O	F
H	156	13	20	77	38	74
O	14	67	12	26	37	30
F	15	2	245	16	7	239

Table 6: Left: Confusion Matrix for predictions of Y-labels of the test set for RF classifier shown in Table 5 (bottom). Right: Confusion Matrix for the same test set but against the true T-labels for the same data points.

The confusion matrices in Table 6 show that predictions produced by the RF classifier trained on the annotated Y-data correspond to a higher degree with the same annotator than the labels between the two annotators. This may indicate that each annotator is consistent in itself, whereas larger disagreement can be observed between the annotators. For instance, in the upper row, out of the 189 Y-annotated H labels, 156 labels were correctly predicted by the trained RF, 13 were predicted as O and 20 as F states. In contrast, annotators Y and T agree in H label only 78 cases. AUs that annotator Y considers H receive in 38 instances label O and 74 cases the label F by annotator T.

4.5 Accuracy across Annotators and Classifiers

In order to corroborate the assumptions in the previous subsection, we assess accuracy patterns across the two classifiers and annotators in more detail. We trained the RF (R) and MLP (M) classifiers with a training set of 1269 AUs to predict T and Y labels, as outlined in section 3. This provided us with four models for the two classifiers (M and R) and two annotators (T and Y). Successively, each of the four models (MT, MY, RT, and RY) predicted a list of state labels for the 544 examples in test set. We thus obtain six lists of state label predictions for the test set: four lists of predictions from the four classifiers (MT, MY, RT, and RY) and in addition the original labels from the manual T and Y annotations. Table 7 shows accuracy scores for the 6×6 pair-wise combinations of these label lists. Since accuracy scores are symmetrical (i.e., $\text{Accuracy}(x, y) == \text{Accuracy}(y, x)$), Table 7 only shows the lower part of the rectangular matrix. Note also that $\text{Accuracy}(x, x) == 1$ and that the triangle below the diagonal adds up to 15 accuracy pairs (cells).

	T	Y	RT	RY	MT
Y	.76	1	—	—	—
RT	.87	.72	1	—	—
RY	.77	.86	.78	1	—
MT	.85	.71	.91	.77	1
MY	.73	.80	.74	.85	.76

Table 7: Accuracy scores for different pairs of Classifiers (R and M) and Annotators (T and Y).

The accuracy scores in Table 7 range between .71 and .91. As discussed in sections 4.1, higher accuracy scores are observed for RF than for MLP and for annotator T as compared to annotator Y. However, contrary to what one might expect, the highest accuracy scores are obtained between predictions of two classifiers trained on the same data, rather than between predictions of a classifier and the data population it was trained on. Thus, Table 7 reveals that:

1. highest accuracy scores are observed when comparing predictions of two different classifiers trained on data of the same annotator. Thus the two comparisons: MT/RT and MY/RX produce among the highest accuracy scores of .91, and .85 respectively. These numbers are marked in **bold** in Table 7

2. high accuracy scores, but not quite as high, are also observed when comparing predictions of a classifier evaluated against the manual annotations of the same annotator that the classifier was trained on. Thus the the four accuracy scores: RT/T, RY/Y, MT/T, and MY/Y produce the second highest accuracy scores of .87, .86, .85, and .80 respectively. This is the case discussed in the context of Table 6 (left).
3. as can be expected, the predictions of two classifiers trained on different annotators provides lower accuracy values as compared to those in item 1. and 2 above. These pairs of HOF state label predictions have the following accuracy values: MY/RT:.74, RY/RT:.78, MT/Ry:.77, and MY/MT:.76.
4. surprisingly, even lower is the accuracy between the two manual annotations T/Y. With a value of .76 it is just slightly higher than the accuracy values of a manual annotation and a different classifier in item 5. This is the case discussed in the context of Table 6 (right).
5. the lowest accuracy scores are observed when comparing predictions of a classifier that was trained on one annotator A but evaluated with a manual annotation of the other annotator B. The the four comparisons: MT/Y, RT/Y, MY/T, and RY/T produce the lowest accuracy scores of .71, .72, .73, .77 respectively. These numbers are marked in *italics* in Table 7

The results are somewhat puzzling. Most surprising is perhaps the finding that the output of the classifiers in item 1. are more consistent (higher accuracy) than the the classifiers in 2 and that accuracy values in 3. are higher than inter-rater accuracy in 4.

Provided that a(ny) classifier generalizes and approximates the inherent structure of the manual annotations, there will be some noise in the generalizations. Under this assumption we expect that accuracy values in 2. should be higher than in 1, since the noise of two classifiers (in item 1) would multiply. Presumably, each of the two classifiers (M and R) would 'infer' their own generalizations which, we would assume, are likely less compatible than the classifier's own generalization about the set of manual annotations which the classifier was trained on (as in item .2). Provided that the manual annotations are consistent, i.e., they are

'gold' data, why then do pairs of classifiers trained on the same data produce higher accuracy values as compared to the gold data?

Why would it be the case that predictions from two different models (R and M) produce more consistent predictions as each of the classifiers evaluated against the test data taken from a population that they were trained on? Provided the test data is correct, how is it possible that, despite their very different nature and implementation, RF and MLP make similar but wrong predictions?

Similar surprising is the observation that T/Y inter-annotator accuracy in item 4. is lower than the predictions of the classifier trained and evaluated on two different annotators in item 3.

Also this outcome suggests that the two classifiers may have inferred similar generalizations that somehow capture similarities between the T and Y training sets, but that do not, however, account correctly for the structure of the test set. This idea is corroborated in the accuracy values reported in item 5. which shows that the worst values are obtained by evaluating a classifier on a manual test set of a different annotator.

The results indicate that an evaluation of the classifier on manually annotated data or on automatically generated test sets may lead to different results. The results may also indicate that the training set is perhaps not sufficiently large to capture the instances that are represented in the test set. However, we take it that our experiments validate the epistemological objectivity of the HOF taxonomy, as the classifier perform well on the task at hand.

5 Discussion and Conclusion

Human translation is a complex cognitive process that involves numerous interacting processes. To understand and analyse these processes, one approach to Translation Process Research (TPR) has been to collect and synchronize behavioral data (keystrokes and gaze data) from translation sessions and to segment the flow of data into various kinds of processing units. Several automatic segmentation approaches have been suggested, but as the labels often lack intuitive understanding it is difficult to interpret the data.

A novel higher-order HOF taxonomy has been proposed (Carl et al., 2024) that segments the data into three phenomenal states in which a translator can be: a state of orientation (O) accounts for the

	Ontology	Epistemology
Subjective	EXISTENCE OF THE SUBJECTIVE <ul style="list-style-type: none"> • Reality as I experienced it (intentions, attitudes, pain, beliefs, desires, etc) • Conscious personal experience 	KNOWLEDGE OF THE SUBJECTIVE <ul style="list-style-type: none"> • Reality as it is judged by me (opinions, preferences, etc) • What “I” know to be the case
	EXISTENCE OF THE OBJECTIVE <ul style="list-style-type: none"> • Reality as it exists: physical, spatial, temporal (mountains, molecules, etc.) • Exists independent of perception 	KNOWLEDGE OF THE OBJECTIVE <ul style="list-style-type: none"> • Reality as ”we” describe it: norms, conventions (money, marriage, etc.) • Assertions “we” make about reality

Table 8: Modes of existence according to Searle.

need of ST information input which is characterized by reading-ahead in the ST. In a flow state (F), translations are fluently produced, and the state of hesitation (H) reflects surprise or uncertainty, which is characterized by regressions, re-fixations and text modifications. Together with the HOF taxonomy, (Carl et al., 2024) specify a decision tree that provided criteria for the annotation process.

A small corpus of behavioral data annotated and released (Carl et al., 2024). The annotated data consists of six English-Spanish translation sessions (approximately 900 words) and 1813 HOF-state annotated Activity Units (AUs, (Carl et al., 2016)). Two annotators annotated the data with HOF labels and — after specifying a decision tree and annotation guidelines — annotators reached a good inter-rater agreement.

Given the novelty of the annotation taxonomy, we investigate how well the HOF annotations can be reproduced. We use machine learning (ML) classifiers to validate the ”epistemological objectivity” (Searle, 2017) of the annotation schema. That is, we deploy a Multi-layer Perceptron and a Random Forests classifier to assess the ”objectivity” of the manual annotations, where high accuracy of the ML classifiers would indicate the validity of the underlying HOF annotations taxonomy.

In his discussion about ”modes of existence”, (Searle, 1998; Searle, 2017) makes a distinction between, on the one hand, subjective and objective ways of understanding and, on the other hand, between the epistemology and the ontology of knowledge and reality (see Table 8). Whereas ontology is a branch of metaphysics that deals with the nature of being, epistemology is the branch of

philosophy concerned with the theory of knowledge.

Ontological subjectivity then refers to the idea that subjective experiences is a form of reality, but there may not be an independent, objective reality beyond these subjective constructions (see Table 8). Epistemological objectivity is the idea that certain aspects of reality can be known objectively, independent of my beliefs, perspectives, or interpretations. Objective knowledge can be discovered or verified through rational inquiry, observation, or evidence, regardless of subjective opinions or interpretations.

Despite the fact that consciousness has a subjective mode of existence—and is thus not directly accessible to scientific inquiry—Searle claims that this does not prevent us from having an epistemological objective science of consciousness. While translators experience subjective states of orientation, hesitation and flow, these states, we assume, can be recovered in the behavioral data and studied under epistemically objective conditions. Norms, regulations or—as in our case the HOF annotation taxonomy—can be understood, deployed and objectively verified within observable TPD in a given context. Our results suggest, however, that there might be a gradual slope between epistemological subjective and epistemological objective modes of existence, rather than a binary one. Table 7 suggests that, despite a well-formulated HOF state decision tree as described in (Carl et al., 2024), a perfect agreement between different annotators may not always be possible⁵. Accuracy values,

⁵Similar findings have been reported in countless translation

such as those in Table 7, may thus provide an index for the degree of epistemological objectivity, where higher accuracy values indicate greater epistemological value of the underlying taxonomy (or norm), and thus allow for higher objectivity while lower accuracy values indicate increased possibilities for epistemological subjectivity. Surprisingly, then, our findings indicate that the two different classifiers (MT/RT and MY/R Y) trained on the same data are able to arrive at higher epistemological objectivity as compared to the two human annotators who follow the same annotation guidelines. It suggests that different classifiers are able to generalize the (training) data in a similar way which, however, deviates from generalizations that our annotators from the annotation guidelines and decision trees.

Acknowledgements

The authors would like to thank Masaru Yamada and Yuxiang Wei for extended discussions about some of the topics addressed in this paper.

References

- [Alves and Vale2009] Alves, Fabio and D. Vale. 2009. Probing the unit of translation in time: Aspects of the design and development of a web application for storing, annotating, and querying translation process data. *Across Languages and Cultures* 10(2), pages 251–273.
- [Carl and Kay2011] Carl, Michael and Martin Kay. 2011. Gazing and typing activities during translation : A comparative study of translation units of professional and student translators. *Meta* 56(4), pages 952–975.
- [Carl and Schaeffer2017] Carl, Michael and Moritz Schaeffer. 2017. Sketch of a noisy channel model for the translation process. In *Empirical Modelling of Translation and Interpreting*, pages 71–116. Berlin: Language Science Press.
- [Carl et al.2016] Carl, Michael, Srinivas Bangalore, and Moritz Schaeffer. 2016. New directions in empirical translation process research. *Heidelberg: Springer International Publishing Switzerland*. doi, 10:978–3.
- [Carl et al.2024] Carl, Michael, Yuxiang Wei, Sheng Lu, Longhui Zou, Takanori Mizowaki, and Masaru Yamada. 2024. Hesitation, orientation, and flow: A taxonomy for deep temporal translation architectures. *Ampersand*, page 100164.
- [Dragsted2010] Dragsted, Barbara. 2010. Coordination of reading and writing processes in translation. In Shreve, Gregory M. and Erik Angelone, editors, *American Translators Association Scholarly Monograph Series*, volume 381, pages 41–62. John Benjamins Publishing Company.
- [Hvelplund2016] Hvelplund, Kristian. 2016. Cognitive efficiency in translation. In Martín, Ricardo Muñoz, editor, *Reembedding Translation Process Research*, pages 149–170. John Benjamins Publishing Company, 9.
- [Jakobsen2011] Jakobsen, Arnt. 2011. Tracking translators’ keystrokes and eye movements with translog. In Alvstad, C., A. Hild, and E. Tiselius, editors, *Methods and Strategies of Process Research*, pages 37–55. John Benjamins Publishing Company.
- [Jakobsen2017] Jakobsen, Arnt Lykke. 2017. Translation process research. In Schwieter, John W. and Aline Ferreira, editors, *The Handbook of Translation and Cognition*. Willey.
- [Jia et al.2019] Jia, Yanfang, Michael Carl, and Xiangling Wang. 2019. How does the post-editing of neural machine translation compare with from-scratch translation? a product and process study. *The Journal of Specialised Translation*, 31:60–86.
- [Kingma and Ba2014] Kingma, Diederik P and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [Klings2001] Klings, H.P. 2001. *Repairing Texts. Empirical Investigations of Machine Translation Post-editing Processes*. Kent, Ohio: Kent State U.P.
- [Kumpulainen2015] Kumpulainen, Minna. 2015. On the operationalisation of ‘pauses’ in translation process research. *The International Journal for Translation & Interpreting Research*, 7:47–58.
- [Königs1987] Königs, Frank. 1987. Was beim Übersetzen passiert: theoretische aspekten, empirische befunde und praktische konsequenzen. *Die neueren Sprachen* 2, 2(86):162–185.
- [Lacruz et al.2014] Lacruz, Isabel, Gregory Shreve, and M. 2014. Pauses and cognitive effort in post-editing. In O’Brien, Sharon, Laura Winther Balling, Michael Carl, Michel Simard, and Lucia Specia, editors, *Post-editing of Machine Translation: Processes and Applications*. Cambridge Scholars Publishing.
- [Malmkjaer1998] Malmkjaer, K. 1998. Unit of translation. *Routledge Encyclopedia of Translation Studies*, page 286–288.
- [Mollo2024] Mollo, Dimitri Coelho. 2024. Ai-as-exploration: Navigating intelligence space.
- [Muñoz and Apfelthaler2022] Muñoz, Ricardo and Matthias Apfelthaler. 2022. A Task Segment Framework to study keylogged translation processes. *Translation and Interpreting*, 14.

- [O’Brien2006] O’Brien, Sharon. 2006. Pauses as indicators of cognitive effort in post-editing machine translation output. *Across Languages and Cultures*, 7(1):1 – 21.
- [Schaeffer and Carl2013] Schaeffer, Moritz and Michael Carl. 2013. Shared representations and the translation process. a recursive model. *Translation and Interpreting Studies* 8 (2), pages 169–190.
- [Schaeffer et al.2016] Schaeffer, Moritz, Michael Carl, Isabel Lacruz, and Akiko Aizawa. 2016. Measuring cognitive translation effort with activity units. In *Proceedings of EAMT*. EAMT.
- [Searle1998] Searle, John. 1998. *Mind, Language, and Society*. Basic Books.
- [Searle2017] Searle, John. 2017. Reflections on free will, language, and political power, talks at google.
- [Vieira2016] Vieira, Nunes. 2016. *Cognitive Effort in Post-Editing of Machine Translation: Evidence from Eye Movements, Subjective Ratings, and Think-Aloud Protocols*. Ph.D. thesis, Newcastle University.
- [Zhong et al.2023] Zhong, Ziqian, Ziming Liu, Max Tegmark, and Jacob Andreas. 2023. The clock and the pizza: Two stories in mechanistic explanation of neural networks.
- [Zou et al.2022a] Zou, Longhui, Michael Carl, and Jia Feng. 2022a. Patterns of attention and quality in english-chinese simultaneous interpreting with text. *International Journal of Chinese and English Translation & Interpreting*.
- [Zou et al.2022b] Zou, Longhui, Michael Carl, Masaru Yamada, and Takanori Mizowaki. 2022b. Workshop on empirical translation process research. In *Proficiency and External Aides: Impact of Translation Brief and Search Conditions on Post-editing Quality*, page np. AMTA.

Appendix

A Features of Classifiers

Both classifiers were trained with a list of 34 features, shown in Table 9. The first 15 features, above the double line, prefixed with “TU_”, are copied from the TU of which the AU is part (see Figure 1). These features thus encode the context of the AU. All ”TU_” features relate to behavioral data, concerning the gaze and the keystroke data, and their duration.

The lower 19 features were extracted from and describe properties of AUs. Similarly, most of the AU features characterize the behavioral data within one AU. However, four of these features are related to properties of the translation product and four features include contextual from surrounding

Feature	Description of feature
$TU_logDurTU$	log-transformed duration of the TU
$TU_WinSwitch$	Number of gaze switches between ST and TT
TU_TrtT	Total reading time on the ST
TU_TrtS	Total reading time on the TT
TU_TrtST	ratio $\log((TrtS + 1)/(TrtT + 1))$
TU_TGset	Intersection of words IDs produced in next TU
$TU_PauseDur$	Ratio of $(Pause+1)/(DurTU+1)$
$TU_ParTrtT$	Duration of concurrent TT reading while typing
$TU_ParTrtS$	Duration of concurrent ST reading while typing
$TU_ParFixT$	#fixations during concurrent TT reading and typing
$TU_ParFixS$	#fixations during concurrent ST reading and typing
$TU_InsDelLog$	ratio of deletions and insertions $\log(Del + 1)/(Ins + 1)$
TU_FixT	Number of fixations on TT
TU_FixS	Number of fixations on ST
$TU_FixDist$	log of max. distance in Y-position of fixations on ST window (in pixel) $\log(FixSpanY + 1)$
$Type$	Type of TU as described in Table 1
$Gram5$	concatenation of AU type with the preceding four AU types
Dur	Duration of the AU
$SGnbr$	#ST words for which translations were produced (concerns AU types T4, T5,T6)
$TGNbr$	#TT words produced (concerns AU types T4, T5,T6)
Ins	#Insertions (concerns AU types T4, T5,T6)
$CrossS$	Average Cross values for ST words produced in AU
$CrossT$	Average Cross values for TT words produced in AU
$ProbSgaze$	Average log probability of source words in GazePath
$ProbTgaze$	Average log probability of target Words in GazePath
$ProbCgaze$	Average log CrossS value in GazePath
$ProbSTCgaze$	Average log of joint ST, TT and CrossS value in GazePath
$HSgaze$	Average entropy of ST words in GazePath
$HTgaze$	Average entropy of TT words in GazePath
$HCgaze$	Average entropy of Cross values in GazePath
$HSTCgaze$	Average entropy of joint ST, TT and Cross in GazePath
$Effort$	sum of log duration for context AUs: T4, T5, T6
$Effect$	sum of log duration for context AUs: T1, T2, T8
$Significance$	$Effect$ minus $Effort$

Table 9: List of features used for Classifier.

AUs information. The features *SGnbr* and *TGnbr* indicate how many source and target words were covered in the AU, while *CrossS* and *CrossT* are measures of the distance / reordering between the source and the target words (Carl et al., 2016). Four of the AU features refer to the nearby context of the AU. *Type* is type of AU (see Table 1) *Gram5* is the concatenation of AU type labels, while *Effort*, *Effect*, and *Significance* take into account Effort/Effect properties of the two surrounding AUs as described in Table 1.

In this study, we define *Effort*, *Effect*, and *Significance* for an AU to depend on the type and the duration of the two preceding two AUs and the next AU. The *Effort* of an AU is computed as the sum of $\log(\text{Dur}(\text{AU}))$ for each context-AU of Type T1, T2 or T8 (no keystroke activity is observed). The *Effect* is computed as the sum of $\log(\text{Dur}(\text{AU}))$ for each context-AU of Type T4, T5 or T6. The *Significance* of an AU is then its *Effect* minus its *Effort*, so that more *significant* AUs are characterized by longer stretches of text production.

B Importance of features in RF Classifier

Table 10 shows the 34 features in their order of importance as obtained during RF training. The list of features is ordered with respect to the importance of the T column (annotator T). The “N” column indexes features according to their importance for T, while the column header “O” provided the rank re-ordering of the importance for the Y data. There is a strong correlation between the two importance vectors of T and Y ($R=0.95$), indicating that slight differences in the annotation of T and Y do not seem to have a large impact on feature importance of the RT classification.

The context of AUs seems to be important for classifying their HOF label. Thus, the 15 TU-inherited features (those preceded by “TU_”) make around 50% (49.18% and 49.94%) of the total importance for T and the Y respectively. Adding to this the importance of the features that account for the external context of the AUs, *Gram5*, *Effort*, *Effect* and *Significance*, increases the importance of context-related features to 72.16% and 73.87% respectively. That is, only 28% and 26% of the HOF state classification is due to AU internal characteristics. Those AU-local features are indicated in bold in Table 10. Also note that the first 11 most important features are all ‘context’ features which make up around 58% in the T set (57% in Y).

N	Feature	T	Y	O
1	<i>TU_PauseDur</i>	0.0829	0.0848	1
2	<i>Significance</i>	0.0770	0.0707	3
3	<i>Effect</i>	0.0709	0.0778	2
4	<i>Effort</i>	0.0569	0.0535	4
5	<i>TU_InsDelLog</i>	0.0513	0.0425	7
6	<i>TU_logDurTU</i>	0.0502	0.0345	10
7	<i>TU_FixS</i>	0.0431	0.0398	9
8	<i>TU_Trts</i>	0.0421	0.0516	5
9	<i>Gram5</i>	0.0345	0.0278	13
10	<i>TU_FixDist</i>	0.0344	0.0271	14
11	<i>TU_TrtsT</i>	0.0332	0.0467	6
12	Dur	0.0272	0.0293	11
13	<i>TU_FixT</i>	0.0271	0.0219	18
14	Ins	0.0268	0.0401	8
15	<i>TU_TrtsT</i>	0.0246	0.0226	17
16	<i>TU_ParTrts</i>	0.0229	0.0279	12
17	CrossS	0.0196	0.0250	15
18	<i>TU_WinSwitch</i>	0.0189	0.0154	24
19	ProbCgaze	0.0183	0.0146	27
20	<i>TU_ParTrtsT</i>	0.0182	0.0234	16
21	<i>TU_TGset</i>	0.0181	0.0156	23
22	Type	0.0181	0.0152	25
23	ProbSgaze	0.0172	0.0134	34
24	<i>TU_ParFixT</i>	0.0166	0.0176	22
25	HTgaze	0.0162	0.0145	29
26	HCgaze	0.0161	0.0151	26
27	HSTCgaze	0.0160	0.0139	32
28	<i>TU_ParFixS</i>	0.0158	0.0204	20
29	HSgaze	0.0154	0.0143	30
30	ProbTgaze	0.0153	0.0136	33
32	TGnbr	0.0151	0.0208	19
31	SGnbr	0.0151	0.0203	21
33	ProbSTCgaze	0.0141	0.0140	31
34	CrossT	0.0110	0.0145	28

Table 10: Importance of features for T and Y annotations. Columns T and Y give the percentage for the respective features. Column N indicates the order of importance for the T annotator while O provides the order for Y annotator.

Perceptions of Educators on MTQA Curriculum and Instruction

João Lucas Cavalheiro Camargo, Sheila Castilho, Joss Moorkens

SALIS/ADAPT Centre

Dublin City University

joao.camargo@adaptcentre.ie

sheila.castilho@dcu.ie, joss.moorkens@dcu.ie

Abstract

This paper reports the results of a survey aimed at identifying and exploring the attitudes and recommendations of machine translation quality assessment (MTQA) educators. Drawing upon elements from the literature on MTQA teaching, the survey explores themes that may pose a challenge or lead to successful implementation of human evaluation, as the literature shows that there has not been enough design and reporting. Results show educators' awareness of the topic, awareness stemming from the recommendations of the literature on MT evaluation, and reports new challenges and issues.

1 Introduction

Academia and industry continuously make efforts to assess the quality of machine translation (MT) systems (Way, 2020), typically using automatic evaluation metrics (AEM) or human evaluation (HE) (Castilho et al., 2018), each approach possessing its own strengths and weaknesses. However, to evaluate an MT system with detailed and actionable results, it is vital to use a balanced approach incorporating HE in the process in conjunction with AEMs (Way, 2020). In particular, the inclusion of HE must be carefully employed so as to not generate hyperbolic reports of the capabilities of MT systems in particular scenarios such as in Hassan et al. (2018).

Some studies have recommended more rigorous HE design principles (Toral et al., 2018; Läubli

et al., 2020) not only to dampen hype, but also to identify systems' weaknesses through an analysis of complex linguistic phenomena (Castilho and Caseli, 2023). While it is not recommended to rely solely on AEM-based evaluations (Moorkens, 2022), the literature shows a common tendency to rely on AEMs without HE (Marie et al., 2021; Rivera-Trigueros, 2022) in the MT community. It is understood that MT use must consider the purpose and value of translations and the expected longevity of the content (Way, 2013), which extends to MT evaluation as well (Doherty et al., 2018). In this manner, risks from MT systems such as grammatical errors or inappropriate words/constructions (Koehn and Knowles, 2017), biases in the output (Prates et al., 2020), which can be dangerous for specific domains such as legal and medical (Vieira et al., 2021), can be prevented with rigorous HE incorporation in MT evaluations. Given these risks and the responsibility of implementing a careful evaluation, complementing automatic with HE is essential to ensure AI technology is safe, beneficial and fair (Dignum, 2020). It can be achieved with ethical behaviours adopted by engineers and technology developers (Moorkens, 2022), which can be further refined with the training of stakeholders themselves (Dignum, 2020).

Thus, this paper focuses on the instructional training of MT quality assessment (MTQA), as part of a doctoral study that intends to create and provide training in HE for Natural Language Processing (NLP) master's students. In this paper, we report results from the qualitative findings of a survey aimed at MTQA educators with both TS and NLP educators. It inquired about the educators' attitudes and recommendations regarding HE in MTQA teaching, exploring where HE can be positioned pedagogically, what HE content should

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

be prioritised, and evaluates the practical considerations that may facilitate or hinder the incorporation of HE into an MTQA curriculum focused on NLP students. The survey explores the following key questions:

1. What are educators' attitudes towards MTQA?
2. What approaches can be taught to foster HE in MTQA?

These findings can inform MTQA trainers and curriculum planners in making informed decisions to foster appropriate HE teaching and deployment, and consequently, its use in MTQA.

2 Related Work

Translation Quality Assessment (TQA) is complex, leading to much debate and different definitions of translation quality, especially in regard to translation technologies, such as MT (Castilho et al., 2018).

MTQA in Translation Studies (TS) curricula has been slowly introduced alongside the use of MT (Korošec, 2011; Dejica-Cartis, 2012) from a curricular standpoint (Doherty and Kenny, 2014) to critically use and assess MT (Rossi, 2017; Moorkens, 2018). Technical aspects of MT also became an element of MTQA teaching, such as building an engine (Farrell and others, 2017), mainly with the intent of empowering trainee translators to understand how the systems work in order to facilitate informed decisions when evaluating the output (Kenny and Doherty, 2014).

Studies have shown that translators in training gain MTQA proficiency through error analysis (Venkatesan, 2018; Looock, 2020), and that translators' ability to identify missing contextual information in MT output and select appropriate language for specific domains is crucial (Núñez, 2019; Bulut, 2019). This mirrors evaluation models used in the industry (Castilho et al., 2018), showcasing academia's efforts to prepare translators.

Accordingly, AEM and other measures of HE have been introduced in the classroom in the TS field. Doherty and Kenny (2014) and Moorkens (2018) introduced adequacy and fluency measures in conjunction with error typologies. Post and Lopez (2014) created a platform on which students could rank MT outputs and generate BLEU scores (Papineni et al., 2002), focusing on the correlation of human judgement with the AEM.

Other platforms were used in classroom settings, such as the *Asiya-Online* toolkit (Giménez and Márquez, 2010), which provided automatic scores, and later, *MutNMT* (Ramírez-Sánchez et al., 2021; Ramírez-Sánchez, 2023) for guided building and evaluation of NMT systems. Krüger (2022) proposed Jupyter notebooks to introduce translators to the technical nature of AEMs while generating different scores such as BLEU, METEOR (Banerjee and Lavie, 2005), chrF3 (Popović, 2015), TER (Snover et al., 2006) and BERTScore (Zhang et al., 2019). Macken et al. (2023) demonstrate a case study of teaching MTQA, by using HE through ranking, adequacy and fluency measures, correlating to AEMs provided by MATEO (Vanroy et al., 2023), a platform that generates BLEU, ChrF, BERTScore, BLEURT (Sellam et al., 2020) and COMET (Rei et al., 2020) scores. These studies demonstrate the effort to introduce different evaluation approaches in the translation classroom, and how important accessible evaluation platforms are for training translators.

The importance of integrating MTQA into curricula is underscored by the concept of MT Literacy (Bowker and Ciro, 2019) which entails understanding the operational mechanisms of MT systems to facilitate their use. Krüger (2022) and Macken et al. (2023) echo the importance of MT literacy in equipping professionals to use and evaluate MT effectively. However, the implementation of training is context-based, the pedagogical guiding principles for MTQA education tend not to be structured.

In the context of NLP education, the few studies that mention MTQA do so only to a minor degree (Alm et al., 2016; Martynova et al., 2018; Artemova et al., 2021). This is due to MT being only one component within the broader spectrum of training, with evaluation assuming a secondary role. However, that does not diminish the importance of evaluation in NLP, as the reasons for its lack of implementation in training may due to absence of space in the curriculum and the lack of structured information on evaluation (Madureira, 2021). As such, organising the insights and recommendations of MTQA educators, both from NLP and TS may lead to fostering MTQA education.

3 Methods

To collect information regarding educators' insights and suggestions on MTQA, an online sur-

vey was designed (approved by the university's Research Ethics Committee, with reference DCU-FHSS-2023-015).

3.1 Design

The survey was created on the platform Qualtrics.¹ It was formulated with closed-ended and open-ended questions, divided in blocks:²

- the plain language statement and consent form³
- 13 questions related to the profile of the educators
- four questions related to opinions and attitudes regarding HE
- 11 questions related to general MTQA pedagogy
- six questions related to recommendations of HE for MTQA

3.2 Participants

The participants are MTQA educators from different fields, such as TS, NLP and Computational Linguistics (CL).⁴

The participants were recruited via: i) social media; ii) email via a curriculum analysis examining universities' postgraduate programmes in Europe and; iii) email collection by examining publications related to MTQA teaching. Note that participants data was anonymised.

4 Data Results and Analysis

As this is an ongoing study, the results reported in this paper are qualitative and small-scale, with the intention of being exploratory, to explore possible relationships and patterns (Cohen et al., 2017). While it is known that smaller samples are not ideal for generalisations (Saldanha and O'Brien, 2014), the qualitative components may inform better the results of the survey as it reaches a larger-scale (McMillan and Schumacher, 2010).

Data was visualised on Qualtrics, which affords analysis of both closed-ended and open-ended questions. For the closed-ended questions, Qualtrics automatically created graphs based on the responses to form variables, and the platform

¹Available on: <https://www.qualtrics.com>

²The full questionnaire can be found in Appendix A.

³This explained the research aims, the ethical aspects and how the data is handled

⁴The distinctions between CL and NLP was made to accommodate possible different curricular nomenclature and personal preferences.

allowed a degree of customisation to change the colour of graphs and combine/separate variables (or groups) as needed. For the open-ended questions, Qualtrics lists the responses by variables (or groups), allowing an interpretive qualitative analysis of the data.

4.1 Participants' Background

Data drawn from 27 participants were analysed.

Q1 - What is your field? Participants could choose multiple fields to accommodate interdisciplinarity among the educators. 18 participants chose 'Translation Studies' as their field of teaching, five participants chose 'Computational Linguistics', seven participants chose 'Natural Language Processing'.

Among the 28 participants, one participant added 'Speech Processing' as their field, one participant added Human-Computer Interaction as their field and another added 'Computer Science' via the 'other' option. While CL and NLP may have often been used interchangeably in research, they represent different streams of research with different emphases, as Tsujii (2011) demonstrates with their experiment. We also acknowledge that the boundaries may not easily be defined (Luz, 2022). Therefore, methodologically we make no distinction between these two groups, and to aid visualisation, the responses from NLP/CL will be organised and reported as a single group, as such, this leads to nine participants in the NLP/CL group.

4.2 Types of MTQA

Participants were asked about the type of evaluation they teach by answering the question:

Q2 - What types of MT evaluation do you teach? As can be seen in Figure 1, the TS group mostly teaches HE, followed by AEMs and semi-automatic evaluation. When prompted in a follow-up question to explain their comments, the TS educators explained their experience:

- One participant notes that MT evaluation is taught to foster MT literacy leading to better use of the systems.
- One participant has PE as the central type of evaluation, while also teaching HE and AEM to a lesser extent.
- One participant focuses on evaluation through PE.

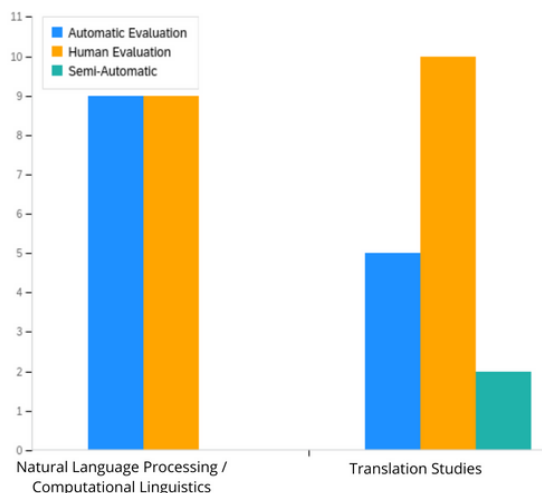


Figure 1: Types of evaluation TS and NLP/CL educators teach

- One participant considers HE the focus of the lesson using the DQF-MQM framework.
- One participant distinguishes MQM from HE, where focus is on MQM, but also mentioning other HE methods, and minor emphasis to AEMs.

Within the NLP/CL group, there are equal efforts reported into teaching HE and AEMs. Four participants described more about their teaching:

- One participant mention teaching HE and AEMs (BLEU, BERT and Comet).
- One participant mentions MTQA is only a component of the course.
- One participant teaches different metrics to different groups. For their Master’s students in Artificial Intelligence, they teach AEMs. For undergraduate translation students, they teach HE.
- One participant mentions teaching AEMs very briefly to make students understand their use in the context of testing the development of a system.

4.3 Attitudes Towards Human Evaluation

This subsection explores participants’ expectations and attitudes towards HE (Q3 and Q4)

Q3 - In your opinion, what trends do you foresee in evaluation metrics that incorporate human judgment for MT systems? Select all that apply.

As may be seen in Figure 2 for the TS group, the most commonly-selected options were context for Quality Assessment (QA), customised evaluation, an equal amount for User Experience (UX) evaluation and multimodal approaches, followed by

ethics, crowdsourced evaluation and two ‘Other’ responses. These two responses were ‘comparing several systems with emphasis on output’ and another response said that all the topics could be important except for crowdsourcing. From NLP/CL, the most commonly-selected were ethics and customised evaluation. Followed by an equal selection of UX evaluation and context-based evaluation. It is worth noting that crowdsourced evaluation was not chosen among the NLP/CL group, which is surprising as the field is known to use crowdworkers for evaluation. The bigger focus given to ethics supports Moorkens’ (2022) assertion that bigger emphasis must be given to the ethical behaviours of engineers, possibly showing that NLP/CL teachers are aware of this. One participant chose ‘Other’ to suggest the use of Large Language Models (LLM) to emulate HE.

Q4 - In your view, what constitutes a comprehensive evaluation of an MT system? Please describe the key components or criteria that should be included.

From TS, nine responses focused only on human judgements and six responses included the use of AEM combined with HE. From NLP/CL, six responses mentioned only human judgements, two responses mentioned a combination of AEM with HE and among the eight answers, four mentioned evaluating MT systems for a specific purpose.⁵ The responses from the TS group mentioned:

- Combined measures of HE and AEM, with State Of The Art (SOTA) metrics, and their correlation.
- Evaluation with platforms with good UX (clean interface, resembling the working environment of a translator).
- Genre, style, terminology, purpose of the text, and agreement with the clients’ needs.
- Use of DQF-MQM for measuring error typology.
- Different degrees of use of MT output, from raw MT to PE at different levels.
- Evaluations that consider human translations as references.
- Measurement of technical aspects (such as training data, speed, pricing, pollution).

The attitudes from the TS group echoes some of the expectations from the industry, such as the adoption of TQA frameworks such as DQF-MQM,

⁵The full qualitative results are included in Appendix B.

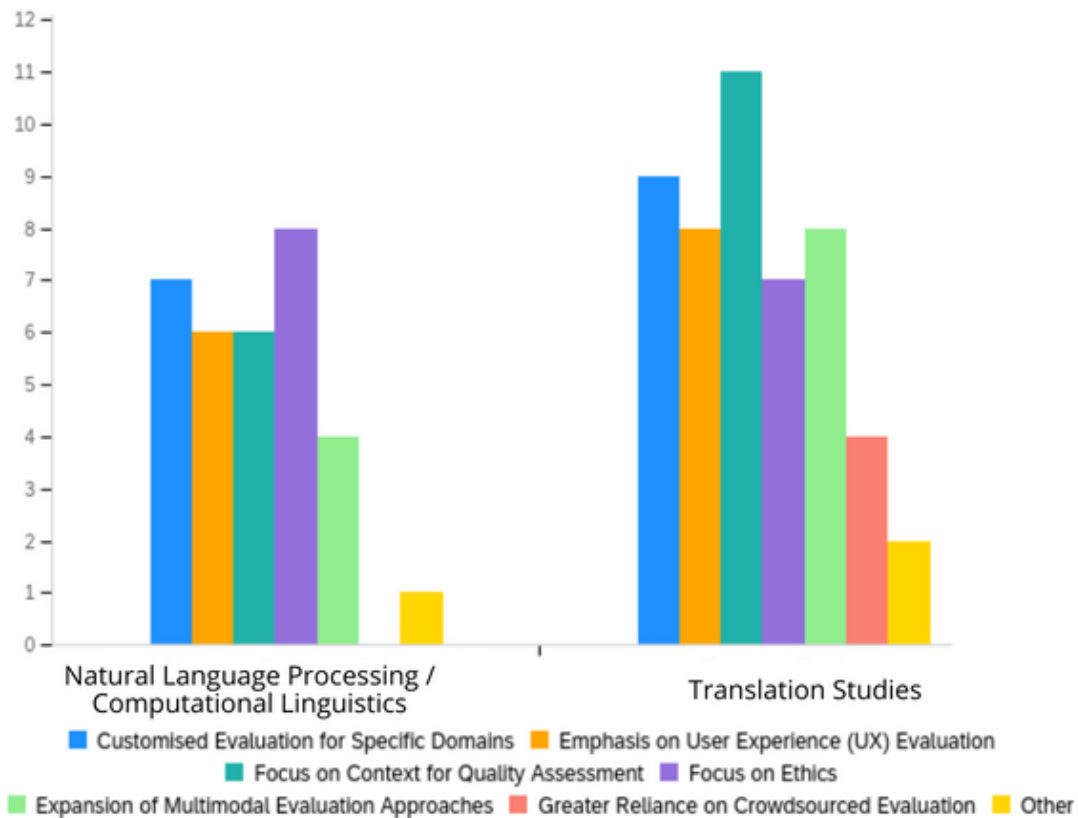


Figure 2: Future trends of Evaluation chosen by TS and NLP/CL teachers

pricing and productivity expectations in addition to the expectation of clients' needs (Castilho et al., 2018). While the NLP/CL group cites the following:

- HE measures such as adequacy, fluency, error analysis and different classifications
- Document-level considerations, such as cohesion and coherence.
- A combination of HE and AEMs, but ultimately with task-based evaluation in mind, to consider how good is the MT system for its appropriate use.
- Risk assessment, considering the type of errors and their severity, according to the domain.
- User-centred assessment, where the end user's purpose of using the translation is to complete a task or is satisfied by its use.

The perishability of content and its purpose (Way, 2013; Way, 2020), in addition to risk assessment which should increasingly be introduced in the training (Doherty et al., 2018) can be noticed by the results of these expected trends. Further, document-level considerations also follow the recommendations made for MT evaluation (Läubli et al., 2020).

4.4 Pedagogical Factors and Recommendations for MTQA

This subsection focuses on the central aspect of MTQA teaching and NLP education (Q5, Q6, and Q7).

Q5 - Assess the importance of including Evaluation Metrics in your academic curriculum - In response to this question, participants assessed the inclusion of both AEM and HE in their teaching curriculum, as can be seen in Figure 3.

Regarding AEMs, the consensus among the NLP/CL participants were that AEMs are 'extremely important', while for TS the most chosen option was 'moderately important'. Regarding HE, while all groups claimed it to be 'extremely important', the TS group mentioned that the emphasis is on HE since they are teaching translators, and therefore AEMs are given less focus. The NLP/CL group mentioned the importance of both AEMs and HE. Interestingly one participant of the TS group mentioned that AEMs are equally important, and one NLP/CL participant stated that, since the course they teach is technical, less emphasis is given to HE.

Following Q5, participants were able to add

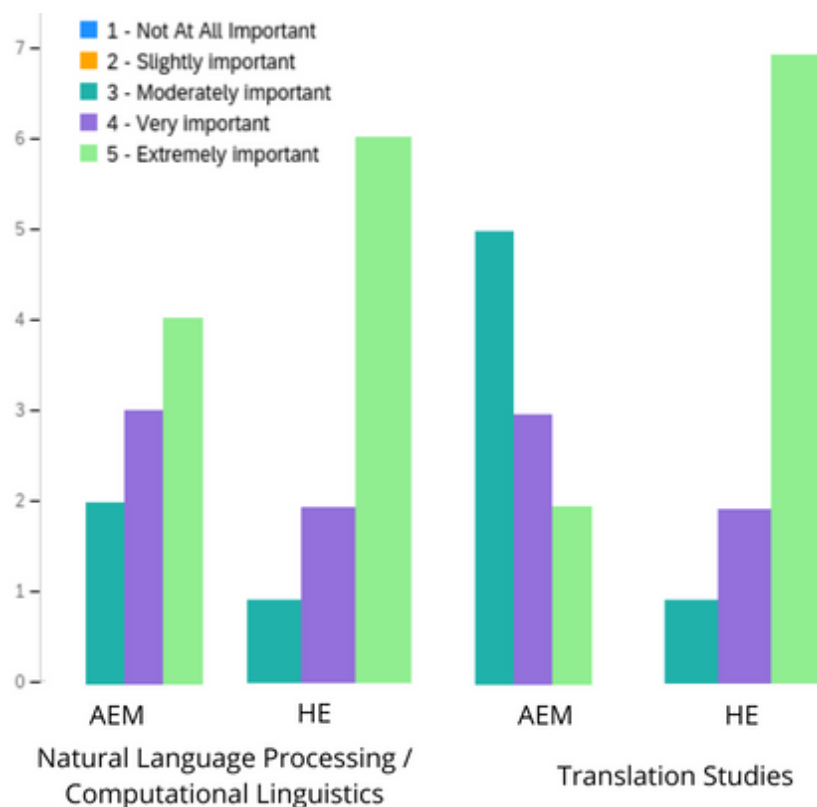


Figure 3: Importance of including AEMs and HE in the curriculum responded by MTQA teachers

comments by responding to *Q6 - Please, add any further comments or explanations for your previous answer*. In the TS group, a participant explained that contextually it is more valuable for them to teach HE towards translators, as AEMs are given less focus. Another participant emphasised that the type of student and level matters when teaching each type of metric. For such participant, undergraduate students who are studying to become translators may require less attention towards both metrics, but the educator explains that for master's NLP students there is room to introduce it to them.

In the NLP/CL group, two participants explained that both metrics are relevant, AEMs providing fast, cheap and objective system checks, while HE are used to understand the values of AEMs and providing insights to improve the systems. One participant differentiates the teaching of metrics in two ways: the first, being moderately important, teaching the metrics directly (such as adequacy scores, error annotation for HE and AEMs such as COMET); while another participant mentions that the most important is to teach the general concepts of HE and AEMs in detail - alluding to a better understanding of the evaluation

process as more important than teaching individual metrics. One participant comments that considering they teach more technical courses, there is less focus on HE.

Both groups correspond to the expectations to a curriculum focused on MT and its evaluation, as what matters the most is the context in which they are inserted (Kenny and Doherty, 2014), whether they are translators or developers, but not forgetting th

Q7 - If you were to create a Human Evaluation module in MT quality assessment addressed to NLP students at Master's level, what should be the main content? Select as many as necessary. In this question participants gave their opinion on the important contents to be taught, as seen in Figure 4.

From the group TS, the most widely chosen option was translators as expert evaluators and design of MT evaluation, followed by adequacy/fluency measures and error typology. The responses from TS may follow the recommendations from the literature such as Laubli et al. (2020) and overall correspond to the importance given to translators (Kenny and Doherty, 2014), such as advisors on the evaluation process (Moorkens, 2017).

From the NLP/CL group, the most widely chosen options were adequacy/fluency measures and inter-annotator agreement, followed by usability and design of MT evaluation. When asked to add other topics (if any), one participant from TS suggested understanding evaluation tools and platforms with analytics, and another TS educator suggested how to interpret results, including generalisability. Among the NLP/CL group, one participant suggested that a whole module on evaluation is not justified, and another participant suggested 'mid-level evaluators', reproducibility of evaluation and bias detection.

Further recommendations in the realm of UX are given, as one of the participants mention how tools and platforms with analytics and insights are important to be introduced, especially if they are accessible. This may be a reason why platforms such as MATEO are being adopted in the classroom (Macken et al., 2023), and to avoid issues that had happened before as reported in Doherty and Kenny (2014) when students were not able to perform AEM scoring due to the unfriendliness of the platforms.

4.5 Pedagogical Challenges

The literature indicates different reasons that may impede more training on evaluation, such as the curriculum (Madureira, 2021) or limited motivation to perform and understand QA processes (Doherty et al., 2018). Thus, this section focuses on the pedagogical elements that may introduce problems in implementing MTQA teaching.

Q8 - Beyond content (such as human evaluation metrics or automatic evaluation metrics), what other pedagogical aspects do you believe may be currently lacking in the teaching of MT quality assessment? Please select all that apply. Participants could select different aspects of teaching such as instructional constraints, hours, and others, as seen in Figure 5

In the TS group, the most commonly-chosen options were allocated hours and faculty expertise and development, followed by curricular structure and lastly by scalability of teaching methods. Expertise and development being one of the most chosen resonates with Doherty et al. (2018) mentioning how educators have to face an evolving and rapidly changing technological scenario, which may make teaching MTQA more difficult. The allocated hours being also one of the most

chosen might be related to MTQA being taught under modules on translation technologies where MT is one component and MTQA is a minor aspect, or a module focused on MT which covers different paradigms, use-cases and MTQA may have more room.

In the NLP/CL group, the most commonly-chosen was allocated hours followed by curricular structure followed by the allocated hours, which has been seen in the literature beforehand as an issue (Madureira, 2021).

Q9 - If you were to create a Human Evaluation module in MT quality assessment addressed to NLP students at Master's level, what would be the best format? Inquired about an ideal format for MTQA training focused on HE, participants responded the following as per Figure

The TS group by a majority suggested an academic module (which is probably unlikely given the previously mentioned time constraints within programmes), followed by the option of a week-long course and a workshop. The NLP/CL group suggested equally an academic module and the option 'other', followed by a two-day course and a week-long course. The 'other' response suggested that each format could be taught depending on the purpose, such as the massive open online course in order to have more time, or a whole-day workshop to introduce the basis of evaluation, or in between a two-day and a week long course, leading the learning to be more contextual. As a follow up, they were asked a question about modality.

Q10 - Given your previous choice on the best format for a Human Evaluation module, what teaching modality would be most suitable? As seen in figure 7, the TS group chose in-person, spread out over several weeks, followed by blended, with the least chosen as an online, synchronous training. Most of the NLP/CL group chose an in-person intensive training, followed by and online synchronous training and an in-person training spread out over several weeks. The 'other' option chosen by a participant of the NLP/CL group suggested that the best modality depends more on the teacher than the topic itself.

Q11 - Please, add any further comments or explanations for your previous answers from Q10 and Q09 here. Within the TS group, one participant commented that the in-person contact is important for the possibility of providing technical

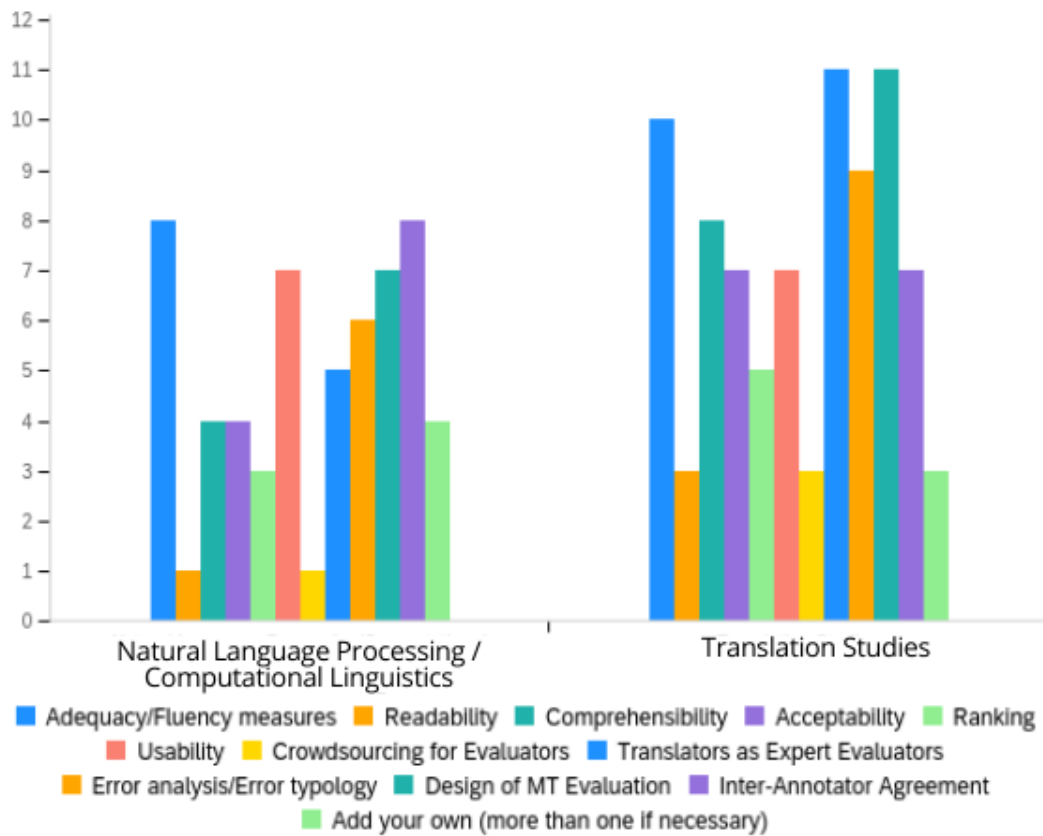


Figure 4: Human Evaluation methods and metrics divided among TS and NLP

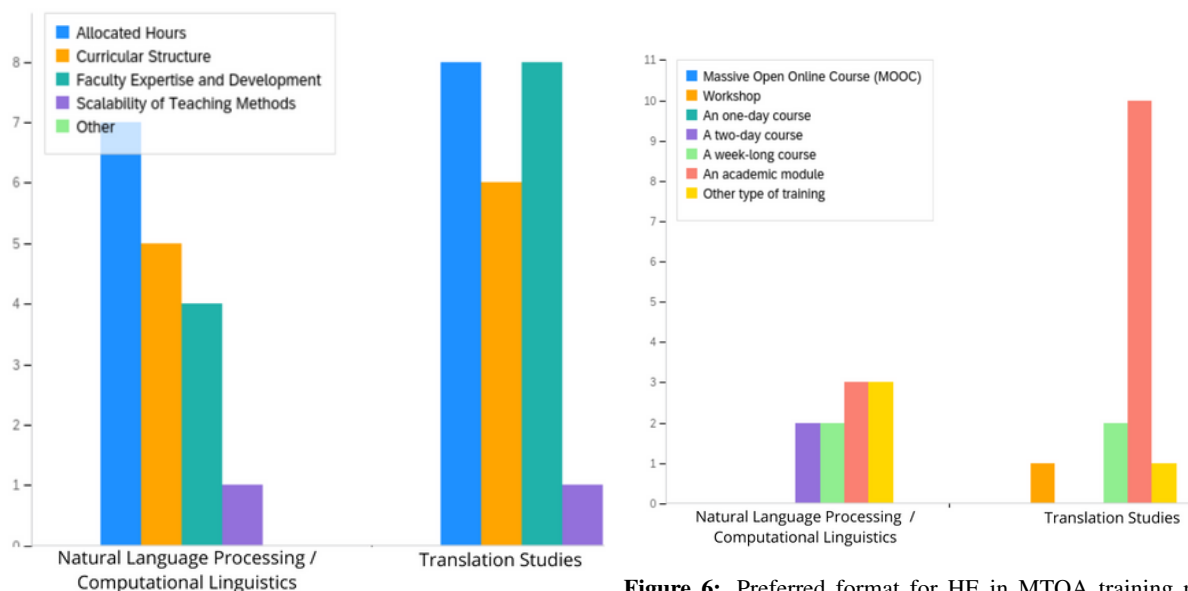


Figure 5: Pedagogical constraints in MTQA teaching divided by TS and NLP/CL

Figure 6: Preferred format for HE in MTQA training responded by TS and NLP/CL educators

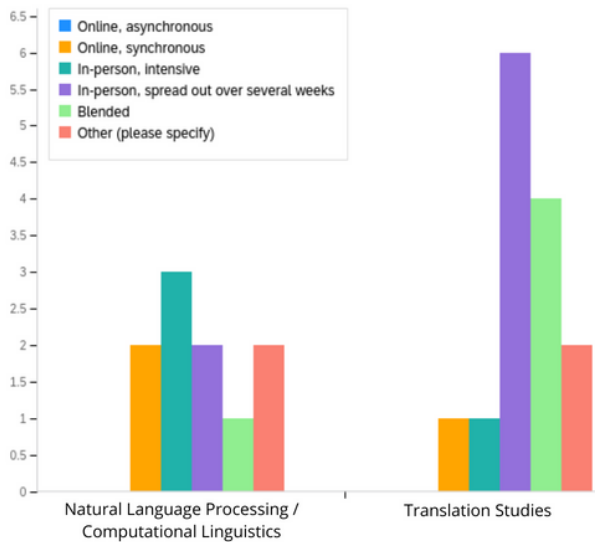


Figure 7: Preferred modality for HE training responded by TS and NLP/CL teachers

support, alluding to easier technical support to students with certain aspects of MTQA. Another participant explained in detail about their experience for a Master’s level training, addressing that academic modules are the only mandatory elements, so the participant suggests spread-out hands-on sessions, such as workshops, in order to provide the different aspects of evaluation for NLP/CL students. Another TS educator complemented that since translator competence takes time to develop, translation evaluation also follows, thus, advocating for long-term training. One educator emphasises that understanding and agreement with the needs of the students would be important to choose the format, so long there was interaction. While the NLP/CL group suggests as many laboratory and hands-on sessions as possible, while another educator suggests that long-term training spanning overall several weeks allow discussion and the opportunity of individual work.

5 Final Considerations and Future Challenges

By comparing the two groups, it can be seen that their attitudes and difficulties reflect both contextual factors of their teaching, and needs commonly associated with their profession.

For TS educators, there has been an increasing effort to integrate the newest technological advancements into their teaching while still maintaining the critical approach of their use. TS educators focus on teaching MTQA for translation

trainees in order to foster their MT literacy, either for more proficient use when performing PE or to prepare them to serve as consultants in the development of MT systems. For either, it places TS educators and the future translators in a position to ensure a safer use of translation systems. NLP/CL educators tend to place more attention towards the ethics, and regard the design of MT evaluation among the most chosen topics, which can be performed by translators who can serve as experts on this process.

We have seen in section 4.5 that the technical aspect may present different pedagogical challenges for TS educators **Q8**, since teaching technical elements to a non-technical audience requires accessible resources. Therefore, there has been research done focused on the experience of translators performing PE, and evaluating MT systems. As a result, over the past years platforms such as MutNMT and MATEO are paramount to make aspects of evaluation accessible, especially when teaching AEMs. Accordingly, for TS lecturers who reported faculty expertise and development as a pedagogical difficulty, those platforms are an important resource for educators.

The NLP/CL group reports other difficulties with MTQA, primarily in finding room in the curricular structure to focus on evaluation (Section (4.5, **Q8**). It is worth noting that this group does not recognise either type of evaluation as less important. In fact, the survey shows that NLP/CL educators recognise the importance of HE in MTQA and teach different evaluation metrics to different groups according to their profiles and roles in the evaluation process. Due to the amount of technical content in development to be covered, it has been suggested by NLP/CL educators that the most appropriate way to cover evaluation would be through intensive, interactive, hands-on workshops to practise different aspects of evaluation - either the design planning, different approaches or the annotation. However, based on the results, NLP/CL educators appear to suggest that NLP master’s students who are choosing to work on MT development and evaluation should know the basic approaches and should still place translators at the centre of the evaluation. These results show the efforts of the MT community at demonstrating the importance of every stakeholder in the MTQA process - from developer to evaluator.

The design and implementation of MTQA still

brings challenges (Section 4.5), but TS, NLP and CL educators report it is essential, whether you are training translators or developers. Challenges to overcome may include:

- Teaching the design of a MT system evaluation is important, but also the user-friendliness of the platform or methodology of evaluators, placing UX as a worthy topic to investigate.
- Finding space in the curriculum for evaluation may be difficult, so a solution proposed is the design, development and implementation of practical workshops around MTQA.
- LLM-based evaluations emulating HE may become more common, and thus, educators need to be prepared to teach NLP professionals the appropriateness of using this approach in evaluation.

As observed in section 4.3, (Q4), the survey also provided some insights on what constitutes a comprehensive evaluation of MT, demonstrating the awareness of the educators.

- Due to its situational nature, the purpose of the system and its end-user are important factors in designing an evaluation.
- A combination of HE and AEM and its correlations are ideal, particularly to show in training.
- Risk assessment and perishability of content are a factor to note the degree of how comprehensive the evaluation should be.

This survey shines light on the directions of MTQA education according to educators from different fields. We hope the insights and recommendations presented here can aid the MT community in fostering MTQA education.

Acknowledgements: We would like to thank the participants for participating in this research. It is a voluntary survey and the results have shown how busy educators are, we present the utmost gratitude and hope the results are able to provide useful pedagogical insights.

Funding: This research was funded by the School of Applied Language and Intercultural Studies at Dublin City University and with the financial support of Science Foundation Ireland at

ADAPT, the SFI Research Centre for AI-Driven Digital Content Technology at Dublin City University [grant number 13/RC/2106_P2].

References

- Alm, Cecilia Ovesdotter, Kathryn Womack, Anne Haake, and Timothy Engström. 2016. A pedagogical model for computational linguistics across curricular boundaries. *Language and Linguistics Compass*, 10(7):335–345.
- Artemova, Ekaterina, Murat Apishev, Veronika Sarkisyan, Sergey Aksenov, Denis Kirjanov, and Oleg Serikov. 2021. Teaching a massive open online course on natural language processing. *arXiv preprint arXiv:2104.12846*.
- Banerjee, Satanjeev and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- Bowker, Lynne and Jairo Buitrago Ciro. 2019. *Machine translation and global research: Towards improved machine translation literacy in the scholarly community*. Emerald Publishing Limited.
- Bulut, Senem ÖNER. 2019. Integrating machine translation into translator training: towards ‘human translator competence’? *transLogos Translation Studies Journal*, 2(2):1–26.
- Castilho, Sheila. and Helena de Medeiros Caseli Caseli. 2023. Tradução automática. *Processamento de Linguagem Natural: Conceitos, Técnicas e Aplicações em Português*.
- Castilho, Sheila, Stephen Doherty, Federico Gaspari, and Joss Moorkens. 2018. Approaches to human and machine translation quality assessment. *Translation quality assessment: From principles to practice*, pages 9–38.
- Cohen, Louis, Lawrence Manion, and Keith Morrison. 2017. *Research Methods in Education*. Routledge.
- Dejica-Cartis, Daniel. 2012. Developing the electronic tools for translators syllabus at politehnica university of timisoara. *Procedia-Social and Behavioral Sciences*, 46:3614–3618.
- Dignum, Virginia. 2020. Responsibility and artificial intelligence. *The oxford handbook of ethics of AI*, 4698:215.
- Doherty, Stephen and Dorothy Kenny. 2014. The design and evaluation of a statistical machine translation syllabus for translation students. *The Interpreter and Translator Trainer*, 8(2):295–315.

- Doherty, Stephen, Joss Moorkens, Federico Gaspari, and Sheila Castilho. 2018. On education and training in translation quality assessment. *Translation quality assessment: From principles to practice*, pages 95–106.
- Farrell, Michael et al. 2017. Building a custom machine translation engine as part of a postgraduate university course: a case study. In *Proceedings of the 39th Conference Translating and the Computer*, pages 35–39.
- Jiménez, Jesús and Lluís Màrquez. 2010. Asiya: An open toolkit for automatic machine translation (meta-) evaluation. *Fifth Machine Translation Marathon*, 94.
- Hassan, Hany, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, et al. 2018. Achieving human parity on automatic chinese to english news translation. *arXiv preprint arXiv:1803.05567*.
- Kenny, Dorothy and Stephen Doherty. 2014. Statistical machine translation in the translation curriculum: overcoming obstacles and empowering translators. *The Interpreter and translator trainer*, 8(2):276–294.
- Koehn, Philipp and Rebecca Knowles. 2017. Six challenges for neural machine translation. *arXiv preprint arXiv:1706.03872*.
- Korošec, Melita Koletnik. 2011. Applicability and challenges of using machine translation in translator training. *ELOPE: English Language Overseas Perspectives and Enquiries*, 8(2):7–18.
- Krüger, Ralph. 2022. Using jupyter notebooks as didactic instruments in translation technology teaching. *The Interpreter and Translator Trainer*, 16(4):503–523.
- Läubli, Samuel, Sheila Castilho, Graham Neubig, Rico Sennrich, Qinlan Shen, and Antonio Toral. 2020. A set of recommendations for assessing human-machine parity in language translation. *Journal of artificial intelligence research*, 67:653–672.
- Loock, Rudy. 2020. No more rage against the machine: how the corpus-based identification of machine-translationese can lead to student empowerment. *The Journal of specialised translation (JoS-Trans)*, 34:150–170.
- Luz, Saturnino. 2022. Computational linguistics and natural language processing. *The Routledge Handbook of Translation and Methodology*, pages 373–391.
- Macken, Lieve, Bram Vanroy, and Arda Tezcan. 2023. Adapting machine translation education to the neural era: A case study of mt quality assessment. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 305–314.
- Madureira, Brielen. 2021. Flamingos and hedgehogs in the croquet-ground: Teaching evaluation of nlp systems for undergraduate students. In *Proceedings of the Fifth Workshop on Teaching NLP*, pages 87–91.
- Marie, Benjamin, Atsushi Fujita, and Raphael Rubino. 2021. Scientific credibility of machine translation research: A meta-evaluation of 769 papers. *arXiv preprint arXiv:2106.15195*, 1(1):7297–7306.
- Martynova, Irina, Lilia Metelkova, Natalia Gordeeva, Larisa Nikitinskaya, Margarita Emelianova, and Alena Trukova. 2018. The programs of computational linguistics graduate in german and us universities. *Visnyk Natsional'noi akademii kerivnykh kadriv kultury i mystetstv*, 1(3).
- McMillan, James H and Sally Schumacher. 2010. *Research in education: Evidence-based inquiry*. pearson.
- Moorkens, Joss. 2017. Under pressure: translation in times of austerity. *Perspectives*, 25(3):464–477.
- Moorkens, Joss. 2018. What to expect from neural machine translation: a practical in-class translation evaluation exercise. *The Interpreter and Translator Trainer*, 12(4):375–387.
- Moorkens, Joss. 2022. Ethics and machine translation. *Machine translation for everyone*, page 121.
- Núñez, Kenneth Jordan. 2019. Análisis de la percepción, la utilidad y la calidad de los sistemas de ta por parte del traductor en formación. *E-Aesla*, 1(5):391–399.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Weijing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Popović, Maja. 2015. chrF: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.
- Post, Matt and Adam Lopez. 2014. The machine translation leaderboard. *Prague Bull. Math. Linguistics*, 102:37–46.
- Prates, Marcelo OR, Pedro H Avelar, and Luís C Lamb. 2020. Assessing gender bias in machine translation: a case study with google translate. *Neural Computing and Applications*, 32:6363–6381.
- Ramírez-Sánchez, Gema, Juan Antonio Pérez-Ortiz, Felipe Sánchez-Martínez, Caroline Rossi, Dorothy Kenny, Riccardo Superbo, Pilar Sánchez-Gijón, and Olga Torres-Hostench. 2021. Multitrainmt: training materials to approach neural machine translation from scratch. In *TRITON 2021 (Translation and Interpreting Technology Online)*.

- Ramírez-Sánchez, Gema. 2023. Mutnmt, an open-source nmt tool for educational purposes. In *Proceedings of the 1st Workshop on Open Community-Driven Machine Translation*.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. *arXiv preprint arXiv:2009.09025*.
- Rivera-Trigueros, Irene. 2022. Machine translation systems and quality assessment: a systematic review. *Language Resources and Evaluation*, 56(2):593–619.
- Rossi, Caroline. 2017. Introducing statistical machine translation in translator training: from uses and perceptions to course design, and back again. *Revista Tradumàtica: tecnologies de la traducció*, 1(15):48.
- Saldanha, Gabriela and Sharon O'Brien. 2014. *Research methodologies in translation studies*. Routledge.
- Sellam, Thibault, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. *arXiv preprint arXiv:2004.04696*.
- Snover, Matthew, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231.
- Toral, Antonio, Sheila Castilho, Ke Hu, and Andy Way. 2018. Attaining the unattainable? reassessing claims of human parity in neural machine translation. *arXiv preprint arXiv:1808.10432*.
- Tsujii, Jun'ichi. 2011. Computational linguistics and natural language processing. In *International Conference on Intelligent Text Processing and Computational Linguistics*, pages 52–67. Springer.
- Vanroy, Bram, Arda Tezcan, and Lieve Macken. 2023. Mateo: Machine translation evaluation online. In *The 24th Annual Conference of The European Association for Machine Translation (EAMT 2023)*, pages 499–500. European Association for Machine Translation (EAMT).
- Venkatesan, Hari. 2018. Teaching translation in the age of neural machine translation. *APLX 2017 at Taipei Tech-Transformation and Development: Language, Culture, Pedagogy and Translation*, pages 39–54.
- Vieira, Lucas Nunes, Minako O'Hagan, and Carol O'Sullivan. 2021. Understanding the societal impacts of machine translation: a critical review of the literature on medical and legal use cases. *Information, Communication & Society*, 24(11):1515–1532.
- Way, Andy. 2013. Emerging use-cases for machine translation. In *Proceedings of Translating and the Computer 35*.
- Way, Andy. 2020. Machine translation: Where are we at today. *The Bloomsbury companion to language industry studies*, 1(1):311–332.
- Zhang, Tianyi, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. BERTscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Appendix A. Full Questionnaire

Questions in bold are the ones selected for this paper.

- Q1 - Name - (Open-ended)
- Q2 - Email - (Open-ended)
- Q3 - List of Countries - (Close-ended)
- Q4 - What is your highest level of education? - (Close-ended)
- **Q5 (In the survey, Q1) - What is your field?** - (Close-ended)
- Q6 - How many hours do you spend teaching per week? Move the slider according to the amount of hours. - (Close-ended)
- Q7 - What are your other main work activities? - (Close-ended)
- Q8 - In your current teaching role, how much influence do you have over the curriculum, including changes to the syllabus and teaching methods? - (Close-ended)
- Q9 - Please, add any further comments or explanations for your previous answer here. - (Open-ended)
- Q10 - What type of students do you work with, mostly? - (Close-ended)
- Q11 - At what academic levels do you currently teach? Please select all that apply. - (Close-ended)
- Q12 - What modality/modalities do you deliver training in? - (Close-ended)
- Q13 - Have you taught MT quality assessment before? - (Close-ended)
- Q14 - Please rate the significance of incorporating human evaluation into the development of MT systems. Rate on a scale of 1 to 5. - (Close-ended)
- Q15 - Please, add any further comments or explanations for your previous answer here. - (Open-ended)
- **Q16 (In the survey, Q3) - In your opinion, what trends do you foresee in evaluation metrics that incorporate human judgment for MT systems? Select all that apply.** - (Close-ended)
- **Q17 (In the survey, Q4) - In your view, what constitutes a comprehensive evaluation of an MT system? Please describe the key components or criteria that should be included.** - (Open-ended)
- **Q18 (In the survey, Q2) - What types of MT evaluation do you teach?** - (Close-ended)
- Q19 - Please, add any further comments or explanations for your previous answer here. - (Open-ended)
- Q20 - Do you teach evaluation for NLP tasks (e.g. summarisation, speech recognition, sentiment analysis) other than MT? - (Close-ended)
- Q21 - Please, add any further comments or explanations for your previous answer here. - (Open-ended)
- Q22 - How many years have you been teaching MT quality assessment? Move the slider according to the amount of years. - (Close-ended)
- Q23 - Assess the importance of teaching students how to plan evaluations for MT systems in your academic curriculum. Please rate the importance of integrating evaluation planning as part of the academic curriculum for MT quality assessment. - (Close-ended)
- Q24 - Please, add any further comments or explanations for your previous answer here. - (Open-ended)
- **Q25 (In the survey, Q5) - Assess the importance of including Evaluation Metrics in your academic curriculum. Please evaluate the importance of integrating evaluation metrics into the academic curriculum for MT quality assessment. You will be presented with two types of evaluation metrics. Rate on a scale of 1 to 5.** - (Close-ended)
- **Q26 (In the survey, Q6) - Please, add any further comments or explanations for your previous answer here.** - (Open-ended)
- **Q27 (In the survey, Q8) - Beyond content (such as human evaluation metrics or automatic evaluation metrics), what other ped-**

agogical aspects do you believe may be currently lacking in the teaching of MT quality assessment? Please select all that apply.

- (Close-ended)

- Q28 - Please, add any further comments or explanations for your previous answer here. - (Open-ended)
- **Q29 (In the survey, Q7) - If you were to create a Human Evaluation module in MT quality assessment addressed to NLP students at Master's level, what should be the main content? Select as many as necessary.** - (Close-ended)
- Q30 - Please, add any further comments or explanations for your previous answer. - (Open-ended)
- **Q31 (In the survey, Q9) - If you were to create a Human Evaluation module in MT quality assessment addressed to NLP students at Master's level, what would be the best format? - (Close-ended)**
- **Q32 (In the survey, Q10) - Given your previous choice on the best format for a Human Evaluation module, what teaching modality would be most suitable? - (Close-ended)**
- **Q33 (In the survey, Q11, in relation to Q9 and Q10) - Please, add any further comments or explanations for your previous answers from Q31 and Q32 here.** - (Open-ended)
- Q34 - Is there anything else you would like to add? - (Open-ended)

Appendix B. Full responses from Q4 - In your view, what constitutes a comprehensive evaluation of an MT system? Please describe the key components or criteria that should be included.

- P1 - adequacy, error annotation and some classification
- P2 - To evaluate an MT system, we should take into account the training data used (quantity and quality) - this includes the pretraining data if the model is based on a pre-trained model-, the size of the model (number of parameters), the memory footprint, the

speed (inference time). The generalization power and particularly the robustness to domain shift should be evaluated.

- P3 - I think that the evaluation of an MT system cannot be detached from the intended purpose. If the MT system is used to generate draft translations the key thing to evaluate is translation productivity. In the MT system is used for gisting, the key thing to evaluate is the ability of the user of the MT system to perform a task after reading the MT output.
- P4 - Human and automatic evaluation. But ultimately, task-based evaluation is most important: how good is the MT for whom in what situation?
- P5 - For assessing the appropriateness of an MT system, I consider that there are different elements worth considering: 1. The domain of use (e.g. medical, legal, etc.) 2. Translation quality (does the MT system provide "good enough" quality for the domain?) 3. The machine translation user experience (MTUX) (Is a translator the one using the MT system? Any other type of MT user? What are the MT needs of this type of user? Undoubtedly, MT needs will vary among different MT users) Once all these elements have been considered and factored in, an informed decision can be taken, whether X system is appropriate or not for a specific use-case and user type
- P6 - Error analysis, Style preservation, Coherence, Document level aspects
- P7 - A comprehensive evaluation of the usefulness (sometimes called "quality") of an MT system should mimic as much as possible the usage scenario and the indicators of usefulness. For instance, if one wants to use MT to increase the productivity of translators, then evaluation should measure productivity in a scenario which is as similar as possible to that in which translators work. Judging "translation quality" through human judgements (usually produced "in vacuo") is clearly inferior to this approach.
- P8 - 'traditional' sentence-level assessment - document-level assessment - user-centered

assessment : does the translation enable readers to complete a task or otherwise 'satisfies' readers? - error analysis: what type of errors we see, what severity they present, and consequently perform a risk assessment, depending on the type of document and the type of errors found

- P9 - * Oriented to particular MT use (assimilation or dissemination, for example) * Blinded in the sense that humans do not know whether they are evaluating other humans or machines to avoid biases * Measuring productivity in case of MT used by professional translators
- P10 - source text as well as output evaluation
- P11 - Accuracy and style
- P12 - Translation quality assesment, i.e. MT vs human output; evaluation of PE effort; consistent terminology, style; error typology (and several other aspects that I am unaware of at this time and/or may arise in the future)
- P13 - The evaluation should take into account accuracy, appropriateness (genre, style, terminology, etc.), general language quality, alignment with clients' needs.
- P14 - Language level. Choice of terminology. Expression of idiolect. Stylistic clarity. Degree of understanding of the sociolect of the translation. Y
- P15 - I think both automatic evaluation and human evaluation are essential. Automatic evaluation should be performed with a sufficiently large sample using one or more SOTA metrics. Human evaluation should be performed in a platform that facilitate scoring with a clean interface and should mimic as much as possible the working environment of a translator.
- P16 - For an evaluation to be comprehensive, it should cover the multiple dimensions involved in the adequacy of the system, from technical aspects (training data, speed, pricing, pollution...) and linguistic (accuracy, fluency, grammaticality, contextual adequacy...) to the user experience (perception, use, ethics...).
- P17 - Evaluation based on both automatic scores and human judgement, as well as investigations into how well they correlate. Comprehensive human evaluation should include error annotation using an error typology such as MQM, ranking tasks and post-editing.
- P18 - Combination of state-of-the-art automatic metrics and human evaluation, including inter-annotator agreement.
- P19 - Beyond the above (usability, context, ethics, multimodality): adequacy metrics, quality-level differentiation, workflow integrability, data transparency
- P20 - Actually, DQF-MQM is a good example of a comprehensive evaluation of MT system.
- P21 - Accuracy and fluency are basic metrics, but the former especially needs to be measured at document level. Appropriate terminology is vital for most domains. Outputs need to be vetted for unwanted bias. Literary and other creative texts require other criteria to be used (e.g. creativity, appropriateness of fictive dialogue, etc.).
- P22 - biases - user experience - no hallucinations - Skopos
- P23 - - The basic fluency and adequacy criteria - Is the information usable for specific contexts. It seems that most evaluation focuses solely on linguistic quality, but it would be important to also evaluate whether raw MT is usable in some situations. For example, is the information patent or law professionals get from raw MT sufficient for them making judgments about the importance and relevance of that information? This is a common and growing use case, but I haven't seen much research that tests its viability

Comparative Quality Assessment of Human and Machine Translation with Best-Worst Scaling

Bettina Hiebl and Dagmar Gromann

University of Vienna, Austria

{bettina.hiebl, dagmar.gromann}@univie.ac.at

Abstract

Translation quality and its assessment are of great importance in the context of human as well as machine translation. Methods range from human annotation and assessment to quality metrics and estimation, where the former are rather time-consuming. Furthermore, assessing translation quality is a subjective process. Best-Worst Scaling (BWS) represents a time-efficient annotation method to obtain subjective preferences, the best and the worst in a given set and their ratings. In this paper, we propose to use BWS for a comparative translation quality assessment of one human and three machine translations to German of the same source text in English. As a result, ten participants with a translation background selected the human translation most frequently and rated it overall as best closely followed by DeepL. Participants showed an overall positive attitude towards this assessment method.

1 Introduction

Human and machine translation quality and their assessment have been of importance in research and industry alike (Harris et al., 2016). Different concepts in the field of translation studies include those focusing on preserving the purpose of the source text in the translation, such as the Skopos theory (Reiss and Vermeer, 1984), on the target text as central point in the analysis of quality as Ammann (1990), or on pragmatic aspects of translation as House (2015).

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

Quality Assessment (QA) approaches in the field of Machine Translation (MT) include QA frameworks for assessment by humans and by machines. Very well-known automated metrics that compare candidate translations to reference translations are, for example, BLEU (Papineni et al., 2002), NIST (Doddington, 2002), and METEOR (Banerjee and Lavie, 2005). MT Quality Estimation (Specia and Shah, 2018) represents a fairly new approach that instead of using reference translations trains machine learning models to predict the output quality of a specific MT system. Human assessment of machine translation consists of human ranking (Macháček and Bojar, 2013), overall assessment (Bojar et al., 2017) or error classification (Popović, 2018) and is generally considered subjective, time-consuming and therefore expensive.

Best-Worst Scaling (BWS) (Louviere and Woodworth, 1990) is an annotation method that addresses these limitations, since it allows for subjective and time-efficient annotations. Annotators are provided with n items in a set and are asked to select the best and the worst items from the set. With a set of four items, this simultaneously leads to a ranking with one clear best item, two that are better than the fourth, and a fourth worst item. BWS has successfully been applied to annotating emotion intensities (Mohammad and Bravo-Marquez, 2017), evaluating stakeholder priorities in health matters (Hollin et al., 2022), assessing consumer preferences in wine attributes (Stanco et al., 2020), among many other application scenarios. BWS has also been applied to assess gender-fair language strategy preferences in translation (Paolucci et al., 2023). However, to the best of our knowledge, comparative translation quality assessment with BWS has not been proposed before.

In this paper, we propose a comparative analysis of five sets of four German translations of the same English source text, one human and three machine translations from Google Translate, DeepL and Microsoft Bing Translator. Ten master students of translation studies or multilingual technologies selected the best and the worst option from the set and rated the best from +4 to 0 and the worst from 0 to -4 in an online survey. This rating provides an overall score for each translation method, but also allows for a more detailed analysis on how high or low each method is assessed. In contrast to ranking, e.g. Bojar et al. (2013), not each translation needs to be annotated with a rank label for each set, but only the best and the worst. Furthermore, the agreement between choices and ratings can be directly assessed without having to calculate an inter-annotator agreement. The translations were selected across domains and consisted of one paragraph from non-fiction books, which required a comparatively low level of domain expertise from participants. In addition, participants were invited to leave comments on each set in free text fields and evaluate the overall method at the end of the survey. The results showed an overall positive attitude towards BWS. Since translation quality assessment by humans in itself is rather subjective, we believe that BWS provides a viable, time-efficient and easy to implement alternative for comparing translations, which can be a comparison of MT systems, of human translations, or, as in this case study, to compare both.

2 Preliminaries

As a basis for the study presented below, we provide an exemplary overview of selected work on MT quality assessment as well as combined assessment of human and machine translation, with no claims regarding completeness. In addition, we will briefly introduce the concept of BWS and typical use cases.

2.1 Translation Quality Evaluation

The evaluation of translation quality has received much attention in translation studies and is a topic that is open for debate. Proposed methods to quality analysis range from source-oriented functionalist approaches (Reiss and Vermeer, 1984; Nord, 1997) to target-text quality analysis, e.g. Ammann (1990), a focus on pragmatic aspects, e.g. House (2015), and analysis based on comprehen-

sibility dimensions (Göpferich, 2008). A common denominator for translation quality in translation studies and machine translation are the concepts of the source text-focused adequacy or accuracy and the target text-focused fluency (Castilho et al., 2018). The Multidimensional Quality Metric (MQM) (Lommel et al., 2014) proposes a framework for translation quality evaluation to be applicable to human and machine translation alike. To this end, a catalogue of known quality issues that can be used as an error typology is presented. Another similar error typology that also considers automation was proposed by Popović (2018). As a mid-way between error classification and overall rating or ranking, Popović (2020) propose to mark all words, phrases, and sentences of a target text that are problematic in terms of comprehensibility and adequacy.

Well-known automated metrics to evaluate machine translation are BLEU (Papineni et al., 2002), NIST (Doddington, 2002), and METEOR (Banerjee and Lavie, 2005), which rely on existing reference translations. In order to avoid having to use reference translations, the idea of Machine Translation Quality Estimation (MTQE) (e.g. Specia and Shah (2018)) was proposed, in which machine learning models are trained to predict the translation quality of MT models. For instance, COMET (Rei et al., 2020) takes the source text into account in training a multilingual MT evaluation model and seeks to assimilate human rankings. Toral and Way (2018) used both BLEU as well as human assessment. In a similar fashion, Webster et al. (2020) compared English to Dutch literary translations by humans and the NMT-based systems Google Translate and DeepL, assessing them using manual annotation as well as different metrics in order to get insights into lexical richness, cohesion, syntactic and stylistic parameters. They found NMT to follow the sentence structure of the source text more closely and that human translation tends to have more lexical richness and local cohesion. Several others (Ortiz-Boix and Matala, 2017; Jia et al., 2019) focus on comparing human translations to post-edited MT output or on the influence of machine translation errors or quality on post-editing effort and performance (Carl and Báez, 2019; Munkova et al., 2021).

Approaches to quality evaluation that are closest to BWS relate to ranking and rating of translations. In the Workshop on Machine Translation (WMT)

starting from 2013 (Macháček and Bojar, 2013; Bojar et al., 2013) five different machine translations with 30 or less words for one source text were humanly ranked relative to each other, allowing for ties in rank. These collected rank labels were then used to assign an overall score to each MT system. In 2017 (Bojar et al., 2017), moved from pair-wise ranking to direct assessment of one machine translation with a reference translation on a 0-100 rating scale by means of crowd-sourcing. The yearly collocated WMT Metrics Shared Task (Freitag et al., 2023) asked professional annotators to label problematic sequences with an MQM error category and severity to be compared with automated evaluation metrics and approaches. With BWS annotators only select the best and the worst translation from a set instead of assigning error categories, rank labels or scores to all translations. Furthermore, there is no need to additionally calculate inter-annotator agreements, since both the annotator’s choices and numbers assigned to the best and the worst options allow for a direct comparison of translations/systems, especially given the negative scores for worst translations. In other words, the method as such is designed to provide a comparative score between translations and annotators as a result across all sets.

2.2 Best-Worst Scaling

BWS(Louviere and Woodworth, 1990; Louviere et al., 2015) was developed by Louviere in the 1980s for measuring a list of objects by dividing them into subsets, which are measured on one or more underlying, latent, subjective scales by selecting the best and worst option of each set (Louviere et al., 2015), allowing for comparative rating. Its underlying concept is random utility theory (RUT) (Thurstone, 1927), which assumes that humans are rational decision-makers trying to maximize utility when making choices (Cascetta, 2009), but acknowledges that the utilities have a random component (Louviere et al., 2015). Originally applied mostly in the field of psychology, BWS has been used in different fields, such as health, agriculture, environment, business, linguistics, transportation, and other fields, within the last two decades (Schuster et al., 2024). In the context of translation, Balducci Paolucci et al. (2023) conducted a case study focusing on gender-fair language in translation from English to German, using BWS and a Likert scale in order to evalu-

ate preferences of specific gender-fair translation strategies. To the best of our knowledge, the proposed study is the first to use BWS for assessing human and machine translation quality in comparison.

3 Method

In order to assess the translation quality of different machine translation systems as well as human translation, a combination of two methods for measuring subjective assessment was used: Best-Worst Scaling (BWS) and a Likert scale. BWS (Louviere et al., 2015) was used to select the subjectively best and worst translation, whereas the Likert scale (Likert, 1932) was used to rate the quality of the selected best and worst translations. These methods were used jointly in order to not only rate whether a specific translation was perceived to be the best or worst of a set, but also how high or low the selected translations are rated.

3.1 Text Selection

Five original English text passages of non-fiction books as well as their officially published human translations were selected for the case study, in order to guarantee having a good quality human translation as well as texts on different, slightly specific, non-fiction topics, which are history, politics, finance, biology, and physics. The selected texts were taken from the following books:

- Set 1, History: Queen of our Time: The Life of Elizabeth II by Robert Hardman (Pan Macmillan, 2022)
- Set 2, Politics: A Promised Land by Barack Obama (Penguin Books, 2020)
- Set 3, Finance: Bitcoin for Dummies by Peter Kent and Tyler Bain (Wiley, 2023)
- Set 4, Biology: Seven and a Half Lessons about the Brain by Lisa Feldman Barrett (HarperCollins, 2020)
- Set 5, Physics: Quantum Physics for Dummies by Steve Holzner (Wiley, 2013)

Each of the English text passages consists of one to three sentences. The length of the original texts ranges from 31 to 41 words and from 197 to 270 characters. The statistics on the words and characters per text and translation are shown in Table 1.

Set	W EN	C EN	W HT	C HT	W GT	C GT	W DL	C DL	W BT	C BT
1	36	197	38	262	37	251	36	252	34	235
2	46	270	47	345	45	326	45	328	44	317
3	31	206	31	241	35	247	35	249	35	247
4	38	217	52	333	37	261	41	265	38	262
5	46	268	47	329	42	282	42	292	45	299

Table 1: Counts of Words (W) and Characters (C) for the source texts (EN) and the translations by a human (HT), Google Translate (GT), DeepL (DL), and Microsoft Bing Translator (BT)

3.2 Translation Selection

As human translation for each of the texts, the published translation of the books was used. With the exception of *Quantum Physics for Dummies*, for each of the books, only one German translation has been published so far, namely *Queen of our Times: Das Leben von Elizabeth II* (Bastei Entertainment, 2022); *Ein verheißenes Land* (Penguin Verlag, 2020); *Bitcoin für Dummies* (Wiley, 2023); and *Siebeneinhalb Lektionen über das Gehirn* (Rowohlt, 2023). For the book with multiple translations the 3rd edition of the book published in 2020, i.e. *Quantenmechanik für Dummies* (Wiley, 2020), was selected.

For the MT examples, the NMT systems Google Translate¹, DeepL Translate², and Microsoft Bing Translator³ were selected due to their wide usage, popularity, free availability and ease of access.

3.3 Participant Selection

Major criteria for participant selection were a background in translation studies and a very good command of the English and German language. These language skills are required because the English source texts were displayed alongside the German translations in this survey. The selected participants are considered expert annotators in comparison to annotators of other ranking or rating methods that were based on crowd-sourcing (Bojar et al., 2013) or language proficiency (Freitag et al., 2023) without necessarily a professional background in translation, however, the participants selected were no domain experts. The target group consisted of experienced master’s students in their last year of studies. Participants who are currently enrolled in a more technical translation master’s program, were expected to have a bachelor’s degree in translation studies.

¹<https://translate.google.com>

²<https://www.deepl.com/translator>

³<https://www.bing.com/translator>

3.4 Survey Design

After an introductory description of the survey and some general demographic questions and questions on the background/education of the participants, the survey also comprised some questions on MT background and use of the participants. The entire survey including the source text and translations can be found in Appendix A. The tool Questionstar⁴ was used for conducting the survey.

Before starting the main part of the survey, participants were shown a short explanation of BWS and an example of how to rate the selected options. The major part of the survey consisted of five sets of each a source text in English and its four translations to German. Participants were asked to rate the best selected option on a scale from 4 (highest score) to 0 (lowest score) and the worst from from 0 (highest score) to -4 (lowest score). Additionally, participants were invited to provide comments on their choices or the text/translations in a free text field.

The four different translations of the texts were arranged in different order in each set. Reordering translation options between sets is necessary for three main reasons: (1) ensure that participants are not inadvertently biased towards selecting specific options due to translation patterns of individual MT systems, (2) ensure that participants make a reflected choice and not randomly select options, e.g. always first as best and second as worst, and (3) make it harder to be biased by trying to single out a specific choice, which in this case study is the one human translation. The second reason is one very commonly applied for these types of surveys to allow researchers to single out participants that simply click through the sets, without taking the survey seriously. As regards the third reason, the reordering makes sure that participants are not biased towards always selecting the one option where the human translation supposedly is as

⁴<https://www.questionstar.de>

best. This reordering was done using all 24 possible different variations of combining four systems (permutations), reordering them using the *RAND* function in Excel and selecting the first 5 instances for Set 1 to 5. The order per set is shown in Table 2. The order was the same for all participants.

In the last section of the questionnaire, participants were asked about the difficulty of selecting the best respectively the worst translation, for their overall opinion on this method for translation quality assessment, for their experience in assessing translation quality, and for any further comments they would like to share.

In order to evaluate the survey design and measure the approximate completion duration, a PhD student of translation studies was invited to pilot the survey. Especially the length of the chosen texts and their translations are an important factor in the design, since cognitive and temporal overload of participants are to be avoided. The pilot study resulted in an estimated duration of 35 minutes and no negative feedback regarding text length, survey length, or clarity of instructions.

3.5 Analysis

The numeric BWS ratings are summed up by translation option across all sets and all participants and divided by the number of times the item was selected and rated. This provides one overall score for all translation options. Furthermore, the number of times an option was selected at all, as best, and as worst are analyzed and presented. While theoretically it could happen that one option is never selected in the entire survey, this is practically unlikely. However, should this be the case, then the option is considered neither the best nor the worst and results in a score of zero. Additionally, all free text comments, demographic data and other answers were analyzed.

4 Results

In this section, the participants' profiles, their BWS ratings for the five sets of translations as well as the corresponding Likert scale ratings will be presented, followed by an analysis of the free text answers and experience with the BWS method. In total, the overall completion time for the entire survey ranged between 20 and 35 minutes.

4.1 Participants

Out of the ten participants, nine identified as female and one as male; 30% are between 18 and 24 years old, 60% are between 24 and 34 years old, and one person is between 35 and 44 years old. All of them had a bachelor's degree as the highest completed degree, 90% in translation studies and 10% in romance studies. Asked to rate their proficiency in English according to the Common European Framework of Reference for Languages (CEFR), seven candidates selected C2, two candidates C1, and one candidate B2. In addition, they had to rate their proficiency in German according to the CEFR, for which eight candidates selected C2, one C1, and one B2. The expected level of German proficiency for this degree program is C1 (CEFR). Therefore, all candidates have a sufficiently high command of English and German and an education related to languages. This is important, since the survey showed the English source texts alongside the German translations. In addition, nine out of ten participants indicated to have some translation experience and the remaining person to have more than 8 years translation experience.

To complement the profiles, the candidates were asked regarding their use of MT tools. Two candidates indicated to use it once a month, seven several times per week, and one person daily. Regarding the purpose of the use of MT, the selected options were privately (5), work other than professional translation (7), work for professional translation (3). For this question, more than one option could be selected. When asked to indicate whether they have a preferred MT system and if so, which one(s), eight participants mentioned DeepL, one person DeepL and Google Translate, and one person Google Translate. The overall satisfaction with MT quality was indicated as very satisfied (2), somewhat satisfied (6), and neither satisfied nor not satisfied (2). The options not very or not at all satisfied were not selected.

4.2 BWS Ratings

From the four translations across five texts, each translation method was selected more than once as best or worst. Table 3 shows the overall averaged and summed ratings and total number of times each translation method was selected. The best summed rating was attributed to the human translation with 33 points from the Likert scale

Set	Option 1	Option 2	Option 3	Option 4
1	BT	HT	DL	GT
2	GT	HT	BT	DL
3	HT	BT	DL	GT
4	DL	HT	BT	GT
5	HT	DL	BT	GT

Table 2: Order of translations per set

across all sets and participants, closely followed by DeepL with 25. The Microsoft Bing Translator and Google Translation were mostly rated negatively, resulting in a score of -20 for the former and -25 for the latter ranked in the overall last place. As can be seen from the detailed BWS rating results in Table 4, this last position and worst result can be attributed to a collective choice as worst translation by all 10 participants and a very negative score on the Likert scale in Set 2 on politics and an overall low selection rate in other sets (see Table 3). Even the overall best option of human translations was assigned a number of negative ratings across sets, but still achieved enough positive selections and ratings in total to result as best option. It is the overall number of times the option was selected as best/worst and scored highly/poorly that finally counts.

The average rating in Table 3 is calculated as the sum of the positive and negative ratings divided by the total number of times the translation method was selected. The average rating for DeepL amounted to 1.04, being slightly higher than the human translation (0.89), while the average ratings for Google Translate (-1.25) and Microsoft Bing Translator (-1.05) were negative. Overall, human translations were selected as the best version 24 times, i.e., in 48% of all cases, whereas the translations by DeepL were selected as best in 30% of all cases. Google Translate and Microsoft Bing Translator clearly lagged behind. While the former was selected more frequently as best and less frequently as worst than the latter, the scores associated with both options still made the former the overall worst option across sets and domains.

The detailed results of the combined BWS & Likert Scale ratings for each participant and each translation output are shown in Table 4. Each of the ten participants rated one of the four presented translations per set as best and one as worst, resulting in a total of 50 selections for each best (rated

from 4 highest score to 0 lowest score) and worst (rated from 0 highest score to -4 lowest score). Positive ratings are highlighted in green, negative ones in red, and “neutral” ones in gray. In addition, the overall results are color-coded according to the source of the translation, i.e. whether it is human translation or produced by Google Translate, DeepL Translate, or Microsoft Bing Translator.

As shown in Table 4, the human translations (HT) are selected as best option in three sets on history, politics, and biology, as among the worst in Set 3 on finance, and as clearly the worst in Set 5 on physics. The translation output of DeepL receives an overall positive evaluation in four out of five sets, were only Set 4 in biology results in a finally negative rating. In Set 4 on biology 70% of the participants selected it as the worst option. As regards Microsoft Bing Translator, it is evident that it was the least selected best or worst option in total with 19 selections, where Google Translate obtained only one more selection with a total of 20. Both obtained very negative ratings in one set, Set 1 for the former and Set 2 for the latter. Interestingly BT is the only option not to be selected at all in one set. It can be seen from these results that it is not only the number of times a translation mode is being selected, but also the exact scores associated with a translation. A translation selected considerably less frequently than the human translation (37 times as opposed to 24 times) can still obtain rather positive results if the individually, per-set attributed scores are overall more positive.

4.3 Ratings & Participant Comments per Set

The full source texts and the translations are provided as part of the survey shown in Appendix A. For each individual set participants had the option to comment on their choices of best and worst as well as their evaluations of the translations in a free-text field. For Set 1 on history, a paragraph describing the role of the Lord Chamberlain in the

	HT	GT	DL	BT
Sum Rating	33	-25	25	-20
Avg. Rating	0.89	-1.25	1.04	-1.05
Sum Rating without Set 3	36	-22	10	-17
Times Selected Best	24 (48%)	7 (14%)	15 (30%)	4 (8%)
Times Selected Worst	13 (26%)	13 (26%)	9 (18%)	15 (30%)
Total Selected	37	20	24	19

Table 3: Average, summed and total BWS rating results

Target Mode	Set 1 (History)				Set 2 (Politics)				Set 3 (Finance)				Set 4 (Biology)				Set 5 (Physics)			
	1.1 BT	1.2 HT	1.3 DL	1.4 GT	2.1 GT	2.2 HT	2.3 BT	2.4 DL	3.1 HT	3.2 BT	3.3 DL	3.4 GT	4.1 DL	4.2 HT	4.3 BT	4.4 GT	5.1 HT	5.2 DL	5.3 BT	5.4 GT
P1	-3		2		-3	3					2	-3		-2		2		1	-2	
P2		2		-1	-3	1			3	-1			-4	3			-1		3	
P3	-2	1			-2	2			-3		1			4	-4		3	-1		
P4	-2	3			-4	4			-3		2		-4	3			-2		2	
P5	-3	3			-4	3			-2	2			-2	3			-3	4		
P6	0			2	-4			3	2	0			0	0			-4			3
P7	-1		3		-4	3			-2	-1	3		-2	3				3		-1
P8	-3	3			-4			4	-1		4		-1		4		-4	4		
P9	-2			3	-4	4			-1		3		-2			3	-2	-2		
P10	-4	3			-4	2			2	-3				-3		2	-1			1
Sum	-20	15	5	4	-36	22	0	7	-3	-3	15	-3	-15	11	0	7	-12	13	3	3
Best	0	6	2	2	0	8	0	2	3	1	6	0	0	6	1	3	1	5	2	2
Worst	9	0	0	1	10	0	0	0	5	4	0	1	7	2	1	0	6	2	1	1

Table 4: Detailed BWS rating results per participant, strategy, and text (HT = Human Translation, GT = Google Translate, DL = DeepL Translate, BT = Microsoft Bing Translator)

Royal Household, was selected. As shown in Table 4, the translation rated as best most often and receiving the best ratings, is the human translation, whereas the translation by Microsoft Bing Translator is chosen as worst most often and receives the worst rating. The trickiest part of the paragraph for translation was the half-sentence in brackets after describing the Lord Chamberlain as a chairman, saying “it has yet to be a woman”. The human translator opted for translating this as “eine Frau konnte sich für dieses Amt noch nicht durchsetzen” (so far no woman has not yet been able to win this office), whereas the translation by Microsoft Bing Translator reads “es ist noch keine Frau” (it is not yet a woman), the one by DeepL “eine Frau hat es noch nicht gegeben” (there has not yet been a woman), and the one by Google Translate “eine Frau ist es bisher noch nicht” (it is not yet a woman). In evaluating the comments, it turned out that this half-sentence was the crucial reason for participants to select HT as the best and BT as the worst option. Other comments reflected on the different translations for “non-executive”, arguing that “nicht geschäftsführend” sounds more natural than “nicht-exekutiv” and the translations for “Royal Household”, with participants expressing differing opinions on translating it as “königlicher Haushalt” or “Königshaus”.

The paragraph selected for Set 2 on politics

is written by Barack Obama describing how his interest in books provided him with knowledge helping him during high school and college. For this set too, the human translation gets the highest overall ratings (22) and is selected as best option and GT as the worst. The most challenging part of this paragraph according to participants’ comments should be at the end of the sentence, when he refers to “bull sessions”, i.e., informal talks/discussions. While all MT systems translate this literally as “Bullensitzungen”, the human translator uses “Diskussionsrunden” (discussion group meetings). Other than that, the most apparent difference between the translation by Google Translate and all other options commented on by most participants is a problem in authenticity and fluency, with changes in the word order contributing grammatical issues, which finally result in its overall selection as worst translation.

For Set 3 on finance, the selected paragraph is a fairly general one about the influence of the launch of Bitcoin on blockchain and cryptocurrency. For this set, the translations produced by Google Translate and Microsoft Bing Translator were exactly the same, so strictly speaking in this set only three different translations were compared. If these two identical translations were counted as one, the summed rating would be -6, once selected as best and 5 times as worst. These

should have clearly been eliminated/changed by the researchers before distributing the survey. Interestingly, the translation by DeepL, which differs only slightly from the ones by the two mentioned above, was rated as the best 6 times. The overall rating for the human translation in this set is negative (-3) and it was chosen as worst option by 5 participants. Most of the participants mentioned deciding for best and worst either due to the first sentence or due to the ending. Four of the six participants who selected DeepL's translation commented that their choice was because they liked how the first sentence was phrased, which is "im Bereich Blockchain und Kryptowährung" (in the field of blockchain and cryptocurrency), which no other version used. Four out of the five participants who chose the human translation as the worst option, commented that they did not like the last sentence of the text, i.e., the combination of "Achtung" (Attention) followed by a comma and ending the sentence with an exclamation mark. All of these make it sound more colloquial in German than the English source, and one person commented that they did not like the addition of the word "wahre" (actual) to the "revolution" in the first sentence of the human translation, as this changes the tone of the sentence. Some participants commented that they would prefer "Achtung" to the wording used by all three MT systems "seien Sie gewarnt" (be warned).

The paragraph in Set 4 on biology compares energy efficiency to a financial budget. The human translation is selected as best translation 6 times. The Google Translate version has been selected as best three times with a total rating of 7. The translation produced by DeepL with a total rating of -15 is selected as worst most often (6 times). One specificity of the results for this set is that P6 chose HT as best and DL as worst, but rated them both with "0", which indicates that for this participant neither of the translations was particularly good or bad. Interestingly, for this set several participants commented that deciding on the best and worst translations is difficult without the context or information on the use case, as the human translation is translated much more loosely than the others. Those who selected the translation by DeepL as worst, commented that it is either not coherent or incomprehensible, hard to read, or sounds artificial. The most challenging part of this paragraph was the second half of the last sen-

tence, i.e. "tracks resources like water, salt, and glucose as you gain and lose them." As one participant phrased it, the wording DeepL used "wie Sie zu- und abnehmen" in German sounds as if referring to gaining and losing weight. The participants who chose the human translation as the best version commented that it sounded most natural, was easy to understand, and translated in a creative and not too literal way. One participant who selected the human translation as worst commented that, depending on the context, this translation could also be the best translation, but that without context they do not perceive it as faithful to the source text enough.

The paragraph selected for Set 5 on physics was concerned with black bodies and the spectrum of light emitted by them. The human translation was clearly rated worst for this set with a total rating of -12 and the translation by DeepL clearly rated best with an overall rating of 13. With this set, several participants who selected the human translation as worst mentioned that they did not like that instead of opting for the literal translation of the field "physics" it was translated as "die Physiker" (the physicists), which is not only less general than the field but also adds a masculine gender in German. It was also criticized that in the human translation the word "sogenannten" (so-called) was added. However, one participant who chose the human translation as worst argued that the line between best and worst translation was very thin in this case. The arguments for selecting the translation by DeepL as the best one were that it translates "physics" literally while being most fluent, adding no additional words, and being most appropriate and straightforward.

4.4 General Comments by Participants

At the end of the questionnaire the participants were asked about how easy/difficult they found selecting and rating the best and worst translations overall. The selected options for the degree of difficulty to select the best translation were somewhat easy (4), neutral (1), somewhat difficult (4) and very difficult (1). Rating the difficulty of the best translation per set was judged to be somewhat easy by two participants and somewhat difficult by eight, so the big majority found it to be more difficult to rate than to select the best option. Selecting the worst translation was found to be very easy (3), somewhat easy (3), neutral (1), somewhat dif-

difficult (2), and very difficult (1). Rating the translation selected as worst was considered to be somewhat easy by four participants, neutral by three, and somewhat difficult by three.

Participants' free-text comments on the comparative approach of BWS for choosing the best and worst translation were mostly positive. Only one participant mentioned that the approach does not allow for assessing the options regarding more than one dimension making the approach much more subjective. Overall, the approach was found to be a promising and interesting approach and useful for translation quality assessment. According to the participants, the availability of different solutions make you aware of several ideas you would not have considered on your own. Also, it was commented that finding the worst option was easiest and finding the best option much harder, but all in all assessment was easier with more options than having to grade one option would be, although sometimes the best option might have been a combination of several of the available options. It was also commented that in some sets context was missing and might have changed the outcome.

Regarding their experience with assessing (human or machine) translation quality, one participant indicated having used the MQM error typology and one participant having used the MQM-DQA metrics before. Three participants indicated that they do translation quality assessment to a certain degree in translation classes. Several indicated to use some sort of quality assessment on MT or comparing different alternatives for translations of specific sentence parts before deciding which one to use. Participants further indicated that translation quality assessment is fairly subjective except when there are indisputable errors, but also stressful and tiring in general.

General comments on the survey included that the difficulty of selecting and rating the best and worst translations differed for the sets, which is why it would have been better to have the questions on the difficulty for each set rather than in general at the end. Two participants commented that they assume that one of the translations was always human, which they attributed to the fact that it was less faithful/close to the original, arguing that more context would be needed for more reliable decision-making. Since the instructions explicitly stated that the comparison was between human and machine translations, it could be that

the one human translation included per set stood out so clearly that this fact became evident to these participants or that this was the expectation by the participants.

5 Discussion

The major objective of the proposed study was to evaluate whether the method of BWS can effectively be applied to a comparative analysis of translation by humans and/or machines. The purpose of this case study was also to show that this method directly leads to a comparative result without the need of any further inter-annotator agreement calculations or scoring methods of the individual participating translation modes. Overall, it can be stated that the results show a very clear preference for human translations and DeepL from the set of selected MT systems. While a particularly unpopular result for a single set can influence the overall rating, the trends for the participating translation modes are still clearly visible from the final overall results. The one extremely negative rating for Google Translate and Microsoft Bing Translator might have contributed to the overall negative rating, however, the fact that both were not selected as often as the other modes contributed just as much. This statement can be particularly reinforced by the fact that also the human translations ranged among the worst for particular sets. Thus, even though individual sets might influence the final result, the overall tendency of being a viable or less preferable translation option can be deduced from the results. The decision to indicate to participants that there is a direct comparison between different translation modes, i.e., human and machine translation, is entirely open for the proposed method and could easily be adapted.

Human translations were overall selected and rated as the best option, however, it should be noted that each set contained a translation by a different professional translator. This can particularly be noticed by the strong differences in ratings across the sets, which, however, could have been the case with the same human translators for each set. Nevertheless, in future work, it would be interesting to repeat the experiment with human translations from a single professional translator or maybe even two human translations in the individual sets, restricting the number of domains to specific sets of expertise.

BWS is considered a perfectly equipped annota-

tion or prioritization method of subjective nature, which means each person can take a subjective decision. Nevertheless, the overall results return tendencies, especially for translation quality, where at times the selection corresponds to a 100%. As indicated by the comments of the participants, a slight variation in wording or a divergence in the selection of just one word can already influence the decision on whether to select a translation option or not and as best or worst. The advantage of BWS is that strong variations in one set still allow for a tendency and trend in the overall results in the end. Without any further context on the topic, participants selected translations that are more faithful to the source text, which in many cases was one of the MT options. As a matter of fact, the lack of context is the most substantial limitation of this case study, limited source and target texts to less than 50 words. Thus, it should be considered for future surveys how BWS can be provided with more context without risking a cognitive and temporal overload of participants. As a method, it still provides a viable alternative to direct assessment with reference translations and ranking methods, especially considering the number of helpful comments left for this case study.

In terms of limitations, it has to be acknowledged that the number and especially scope of evaluated source texts and translations is strongly limited. Only five sets of individual paragraphs were evaluated in this study. In addition, only the language combination of translating from English to German was considered, which is in favor of training settings of major MT systems. Furthermore, the number of participants was limited to 10. While this is a small number, it, however, shows that BWS is adequate for different sizes of participation numbers. In this study, the objective was less to reach a wide audience but rather to make sure participants have a translation background and experience with quality assessment, in order to test this novel method and to obtain feedback on its efficacy and user-friendliness. In this regard, it is within human nature that it is easier to exclude an option we clearly like least, i.e., select the worst option, rather than identifying the best among four options, which is indicated by the ratings and comments of the participants at the end of the survey. While in this study students with translation backgrounds participated, it would be interesting to repeat this study with professional

and experienced translators, in which case the domain should be limited to their respective expertise. Nevertheless, in this case study, the level of required expertise and technical vocabulary was intentionally kept at a low level to facilitate participation by language rather than domain experts.

6 Conclusion

The study showed that assessing translation quality with BWS is an easy to implement and understand method, which can be successfully administered without lengthy explanations and returns interesting results. The two major benefits to be expected from BWS for translation quality assessment are time efficiency and subjective decisions. Even though the selected number and sizes of translations was small, the survey also only took between 20 and 35 minutes to be completed. While each participant made individual choices for each set, the subjective decisions still provide an overall tendency on which translation method and origin might be preferable for these domains and text genres. It is interesting that in the comments participants remarked on the fact that this is a highly subjective exercise, which, however, when evaluating the overall results is not negative. Quite to the contrary, with BWS and the rating of each selected best and worst option the results show that an effective and consistent comparison of translation quality can be achieved with this method.

For future research we suggest using longer texts in order to provide more context for the MT systems, as well as to perform studies on a larger scale and with professional translators of more experience. In addition, it would be interesting to directly compare this quality assessment method with previously, state-of-the-art methods related to ranking and rating the overall quality of translation, be it machine and/or human, which is part of our future endeavours. In addition, it is interesting to see how this method can be applied to different types of methods related to machine and/or human translation, such as pre-editing, post-editing, and specific translation strategies.

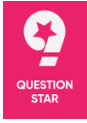
References

- Ammann, Margret. 1990. Anmerkungen zu einer theorie der Übersetzungskritik und ihrer praktischen Anwendung. *TEXTconTEXT*, 5:209–250.
- Balducci Paolucci, Angela, Manuel Lardelli, and Dagmar Gromann. 2023. Gender-fair language in trans-

- lation: A case study. In *Proceedings of the First Workshop on Gender-Inclusive Translation Technologies*, pages 13–23.
- Banerjee, Satanjeev and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan.
- Bojar, Ondřej, Christian Buck, Chris Callison-Burch, Christian Federmann, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut, and Lucia Specia. 2013. Findings of the 2013 Workshop on Statistical Machine Translation. In Bojar, Ondřej, Christian Buck, Chris Callison-Burch, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, Herve Saint-Amand, Radu Soricut, and Lucia Specia, editors, *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 1–44, Sofia, Bulgaria, August. Association for Computational Linguistics.
- Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Raphael Rubino, Lucia Specia, and Marco Turchi. 2017. Findings of the 2017 conference on machine translation (WMT17). In Bojar, Ondřej, Christian Buck, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, and Julia Kreutzer, editors, *Proceedings of the Second Conference on Machine Translation*, pages 169–214, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Carl, Michael and M Cristina Toledo Báez. 2019. Machine translation errors and the translation process: a study across different languages. *Journal of Specialised Translation*, 31:107–132.
- Cascetta, Ennio, 2009. *Random Utility Theory*, pages 89–167. Springer US, Boston, MA.
- Castilho, Sheila, Stephen Doherty, Federico Gaspari, and Joss Moorkens. 2018. Approaches to human and machine translation quality assessment. In *Translation Quality Assessment: From Principles to Practice*, volume 1 of *Machine Translation: Technologies and Applications*, pages 9–38. Springer, Cham, Switzerland.
- Doddington, George. 2002. Automatic evaluation of machine translation quality using n-gram co-occurrence statistics. In *Proceedings of the Second International Conference on Human Language Technology Research*, pages 138–145, San Francisco, CA, USA.
- Freitag, Markus, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore, December. Association for Computational Linguistics.
- Göpferich, Susanne. 2008. *Textproduktion im Zeitalter der Globalisierung: Entwicklung einer Didaktik des Wissenstransfers*. Studien zur Translation ; 15. Stauffenburg, Tübingen, 3. Aufl. edition.
- Harris, Kim, Aljoscha Burchardt, Georg Rehm, and Lucia Specia. 2016. Technology landscape for quality evaluation: Combining the needs of research and industry. In *LREC Workshop on Translation Evaluation*, pages 50–54.
- Hollin, Ilene L, Jonathan Paskett, Anne LR Schuster, Norah L Crossnohere, and John FP Bridges. 2022. Best–worst scaling and the prioritization of objects in health: a systematic review. *Pharmacoeconomics*, 40(9):883–899.
- House, Juliane. 2015. *Translation quality assessment: Past and present*. Routledge.
- Jia, Yanfang, Michael Carl, and Xiangling Wang. 2019. How does the post-editing of neural machine translation compare with from-scratch translation? a product and process study. *The Journal of Specialised Translation*, 31(1):60–86.
- Likert, Rensis. 1932. A technique for the measurement of attitudes. *Archives of psychology*.
- Lommel, Arle, Hans Uszkoreit, and Aljoscha Burchardt. 2014. Multidimensional quality metrics (MQM): A framework for declaring and describing translation quality metrics. *Tradumàtica*, 12:0455–463.
- Louviere, Jordan J and George G Woodworth. 1990. Best worst scaling: A model for largest difference judgments [working paper]. *Faculty of Business*.
- Louviere, Jordan J, Terry N Flynn, and Anthony Alfred John Marley. 2015. *Best-worst scaling: Theory, methods and applications*. Cambridge University Press.
- Macháček, Matouš and Ondřej Bojar. 2013. Results of the WMT13 metrics shared task. In Bojar, Ondřej, Christian Buck, Chris Callison-Burch, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, Herve Saint-Amand, Radu Soricut, and Lucia Specia, editors, *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 45–51, Sofia, Bulgaria, August. Association for Computational Linguistics.

- Mohammad, Saif and Felipe Bravo-Marquez. 2017. Emotion intensities in tweets. In Ide, Nancy, Aurélie Herbelot, and Lluís Màrquez, editors, *Proceedings of the 6th Joint Conference on Lexical and Computational Semantics (*SEM 2017)*, pages 65–77, Vancouver, Canada, August. Association for Computational Linguistics.
- Munkova, Dasa, Michal Munk, Katarina Welnit-zova, and Johanna Jakabovicova. 2021. Product and process analysis of machine translation into the inflectional language. *SAGE Open*, 11(4):21582440211054501.
- Nord, Christiane. 1997. *Translating as a purposeful activity: Functionalist approaches explained*. Routledge.
- Ortiz-Boix, Carla and Anna Matamala. 2017. Assessing the quality of post-edited wildlife documentaries. *Perspectives*, 25(4):571–593.
- Paolucci, Angela Balducci, Manuel Lardelli, and Dagmar Gromann. 2023. Gender-fair language in translation: A case study. In Vanmassenhove, Eva, Beatrice Savoldi, Luisa Bentivogli, Joke Daems, and Janiça Hackenbuchner, editors, *Proceedings of the First Workshop on Gender-Inclusive Translation Technologies*, pages 13–23, Tampere, Finland, June. European Association for Machine Translation.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Weijing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA.
- Popović, Maja. 2018. Error classification and analysis for machine translation quality assessment. In Moorkens, Joss, Sheila Castilho, Federico Gaspari, and Stephen Doherty, editors, *Translation Quality Assessment: From Principles to Practice*, pages 129–158. Springer International Publishing, Cham.
- Popović, Maja. 2020. Informative manual evaluation of machine translation output. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5059–5069.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In Webber, Bonnie, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.
- Reiss, Katharina and Hans J. Vermeer. 1984. *Grundlegung einer allgemeinen Translationstheorie*, volume 147. Max Niemeyer Verlag, Tübingen.
- Schuster, Anne L.R., Norah L. Crossnohere, Nicola B. Campoamor, Ilene L. Hollin, and John F.P. Bridges. 2024. The rise of best-worst scaling for prioritization: A transdisciplinary literature review. *Journal of Choice Modelling*, 50:100466.
- Specia, Lucia and Kashif Shah. 2018. machine translation quality estimation: applications and future perspectives. In *Translation Quality Assessment: From Principles to Practice*, volume 1 of *Machine Translation: Technologies and Applications*, pages 201–235. Springer, Cham, Switzerland.
- Stanco, Marcello, Marco Lerro, and Giuseppe Marotta. 2020. Consumers’ preferences for wine attributes: A best-worst scaling analysis. *Sustainability*, 12(7):2819.
- Thurstone, Louis L. 1927. A law of comparative judgment. *Psychological review*, 34(4):273–286.
- Toral, Antonio and Andy Way, 2018. *What level of quality can neural machine translation attain on literary text?*, pages 263–287. Springer.
- Webster, Rebecca, Margot Fonteyne, Arda Tezcan, Lieve Macken, and Joke Daems. 2020. Gutenberg goes neural: Comparing features of dutch human translations with raw neural machine translation outputs in a corpus of english literary classics. *Informatics*, 7(3):32.

Appendix A. Questionnaire



Best-Worst-Scaling NMT vs HT

Q1

Dear participant,

Thank you for participating in this study on the comparison of machine and human translations. The aim of this survey is to assess translation quality of translations produced by machine translation systems and human translators.

For each example, you will be shown the English source text as well as four different translations. In a first step, please choose which of the four options, according to your opinion, is the best translation and which one you consider the worst option. You will then be asked to rate how good the best option is from 4 (highest score) to 0 (lowest score) and how bad the worst option is in your mind from 0 (highest score) to -4 (lowest score). Please also feel free to comment on your choice, e.g. why one translation is better or worse than the others in your opinion.

Completing the survey should take approximately 20 to 25 minutes and is entirely anonymous. Please provide free text comments in English. Your data and the answers you provide will be recorded for scientific purposes only and analyzed and published anonymously.

Thank you for participating and we hope you enjoy it!

Q2

How old are you?

- 18-24 years
- 25-34 years
- 35-44 years
- 45-54 years
- 55-64 years
- Older than 65 years

Q3

What gender do you identify as?

- Female
- Male
- Diverse/Non-binary
- Prefer not to say

Q4



Best-Worst-Scaling NMT vs HT

What is the highest degree you have completed?

- Bachelor's degree (BA, BSc,...)
- Master's degree (MA, MSc,...)
- Diploma (Mag.)
- Doctorate (Dr., PhD)
- Other, please specify as a comment: _____

Q5

Please indicate the field in which you obtained that degree.

Q6

If you have a degree in translation studies/transcultural communication, please indicate your language combination.

A language _____

B language _____

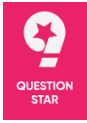
C language(s) _____

Q7

Please indicate your proficiency in **English** according to the Common European Framework of Reference for Languages (CEFR).

- C2 - Proficient user (Mastery)
- C1 - Proficient user (Effective operational proficiency)
- B2 - Independent user (Vantage)
- B1 - Independent user (Threshold)
- A2 - Basic user (Waystage)

Q8



Best-Worst-Scaling NMT vs HT

Please indicate your proficiency in **German** according to the Common European Framework of Reference for Languages (CEFR).

- C2 - Proficient user (Mastery)
- C1 - Proficient user (Effective operational proficiency)
- B2 - Independent user (Vantage)
- B1 - Independent user (Threshold)
- A2 - Basic user (Waystage)

Q9

How would you rate the amount of experience you have translating?

- I have a lot of translation experience
- I have some translation experience
- I have no translation experience

Q10

What is your current profession?

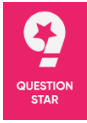
Q11

How many years of professional experience in the translation sector do you have?

- None
- Up to 1 year
- 1 to 3 years
- 4 to 8 years
- More than 8 years

Q12

How often do you use machine translation tools (e.g. DeepL, Google Translate, Microsoft Bing Translator, etc.)?



Best-Worst-Scaling NMT vs HT

- Never
- Up to 5 times per year
- Once a month
- Once per week
- Several times per week
- Daily

Q13

What do you use machine translation for? (please select all that apply)

- Professional translation
- Other work purposes
- Privately

Q14

If you use MT, is there an MT system you prefer? If so, please indicate below which ones you prefer.

Q15

How satisfied are you with the machine translation results in general?

- Very satisfied
- Somewhat satisfied
- Neither satisfied nor not satisfied
- Not very satisfied
- Not at all satisfied

Q16

In each of the following five sections, you will be shown a set of four German translations for one English paragraph. Please indicate which of the four translations - according to your opinion - is the **best** translation and which one the **worst**. You will



Best-Worst-Scaling NMT vs HT

then be asked to rate the chosen options as well as to motivate your choice. **Please do not select the same translation for best and worst option and make sure to only select one option as best and exactly one as worst, which means two options will remain unrated.** Please find below an example of how it should be done (if there is one box selected in each of the columns, this is correct).

Source text in English	Best	Worst
German translation 1	<input checked="" type="radio"/>	<input type="radio"/>
German translation 2	<input type="radio"/>	<input type="radio"/>
German translation 3	<input type="radio"/>	<input checked="" type="radio"/>
German translation 4	<input type="radio"/>	<input type="radio"/>

Q17

Set 1:

English original: At the top of the Royal Household is the Lord Chamberlain, often likened to a non-executive chairman (it has yet to be a woman). He is appointed on a part-time basis to oversee the whole operation.

	Best	Worst
An der Spitze des königlichen Haushalts steht der Lord Chamberlain, der oft mit einem nicht-exekutiven Vorsitzenden verglichen wird (es ist noch keine Frau). Er wird auf Teilzeitbasis ernannt, um den gesamten Betrieb zu beaufsichtigen.	<input type="checkbox"/>	<input type="checkbox"/>
An der Spitze des Britischen Hofes steht der Lord Chamberlain, oft verglichen mit einem	<input type="checkbox"/>	<input type="checkbox"/>



Best-Worst-Scaling NMT vs HT

	Best	Worst
nicht geschäftsführenden Vorsitzenden (eine Frau konnte sich für dieses Amt noch nicht durchsetzen). Er wird auf Teilzeitbasis eingestellt, um den gesamten Betrieb zu leiten.		
An der Spitze des Königshauses steht der Lord Chamberlain, der oft mit einem nicht geschäftsführenden Vorstandsvorsitzenden verglichen wird (eine Frau hat es noch nicht gegeben). Er wird auf Teilzeitbasis ernannt, um den gesamten Betrieb zu überwachen.	<input type="checkbox"/>	<input type="checkbox"/>
An der Spitze des königlichen Haushalts steht der Lord Chamberlain, der oft mit einem nicht geschäftsführenden Vorsitzenden verglichen wird (eine Frau ist es bisher noch nicht). Er wird auf Teilzeitbasis ernannt, um den gesamten Betrieb zu überwachen.	<input type="checkbox"/>	<input type="checkbox"/>

Q18

On a scale from **0 (lowest score)** to **4 (highest score)** , how good would you rate the translation you selected to be the **best** one?

- 4 (highest score)
- 3
- 2
- 1
- 0 (lowest score)



Best-Worst-Scaling NMT vs HT

Q19

On a scale from **0 (highest score)** to **-4 (lowest score)**, how bad would you rate the translation you selected to be the **worst** one?

- 0 (highest score)
- 1
- 2
- 3
- 4 (lowest score)

Q20

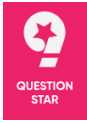
Please comment on what makes the one translation the best and the other the worst in your opinion. Feel free to add any other comments you would like to share on this text or the translations.

Q21

Set 2:

English original: My interest in books probably explains why I not only survived high school but arrived at Occidental College in 1979 with a thin but passable knowledge of political issues and a series of half-baked opinions that I'd toss out during late-night bull sessions in the dorm.

	Best	Worst
Mein Interesse an Büchern erklärt wahrscheinlich, warum ich nicht nur die Highschool überlebte, sondern 1979 auch mit einem dürftigen, aber passablen Wissen über	<input type="checkbox"/>	<input type="checkbox"/>



Best-Worst-Scaling NMT vs HT

	Best	Worst
politische Themen und einer Reihe unausgegorener Meinungen, die ich während der nächtlichen Bullensitzungen von mir gab, am Occidental College ankam der Schlafraum.		
Mein Interesse an Büchern erklärt vermutlich, warum ich nicht nur die Highschool überstand, sondern 1979 beim Eintritt ins Occidental College über ein zwar dünnes, aber einigermaßen passables Politikwissen verfügte und ein paar halb gare Ansichten entwickelt hatte, die ich bei nächtlichen Diskussionsrunden im Studentenwohnheim zum Besten gab.	<input type="checkbox"/>	<input type="checkbox"/>
Mein Interesse an Büchern erklärt wahrscheinlich, warum ich nicht nur die High School überlebte, sondern 1979 auch mit einem dünnen, aber passablen Wissen über politische Themen und einer Reihe von unausgegorenen Meinungen am Occidental College ankam, die ich während nächtlicher Bullensitzungen im Wohnheim ausstieß.	<input type="checkbox"/>	<input type="checkbox"/>
Mein Interesse an Büchern erklärt wahrscheinlich, warum ich nicht nur die Highschool überlebte,	<input type="checkbox"/>	<input type="checkbox"/>



Best-Worst-Scaling NMT vs HT

	Best	Worst
sondern 1979 am Occidental College ankam, mit einem dünnen, aber passablen Wissen über politische Themen und einer Reihe halbfertiger Meinungen, die ich während der nächtlichen Bullensitzungen im Studentenwohnheim in die Runde warf.		

Q22

On a scale from **0 (lowest score)** to **4 (highest score)** , how good would you rate the translation you selected to be the **best** one?

- 4 (highest score)
- 3
- 2
- 1
- 0 (lowest score)

Q23

On a scale from **0 (highest score)** to **-4 (lowest score)** , how bad would you rate the translation you selected to be the **worst** one?

- 0 (highest score)
- 1
- 2
- 3
- 4 (lowest score)

Q24



Best-Worst-Scaling NMT vs HT

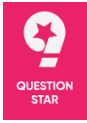
Please comment on what makes the one translation the best and the other the worst in your opinion. Feel free to add any other comments you would like to share on this text or the translations.

Q25

Set 3:

English original: The launch of Bitcoin set off a revolution in blockchain and cryptocurrency. There are now more than 13,000 different cryptocurrencies. (Most, be warned, are essentially valueless and will remain that way.)

	Best	Worst
Der Start des Bitcoin-Netzwerks löste eine wahre Blockchain- und Kryptowährungsrevolution aus. Inzwischen gibt es über 13.000 verschiedene Kryptowährungen. (Achtung, die meisten davon sind im Wesentlichen wertlos und werden es auch bleiben!)	<input type="checkbox"/>	<input type="checkbox"/>
Die Einführung von Bitcoin löste eine Revolution in der Blockchain und Kryptowährung aus. Mittlerweile gibt es mehr als 13.000 verschiedene Kryptowährungen. (Die meisten, seien Sie gewarnt, sind im Wesentlichen wertlos und werden es auch bleiben.)	<input type="checkbox"/>	<input type="checkbox"/>
Die Einführung von Bitcoin löste eine Revolution im Bereich Blockchain und Kryptowährung aus. Inzwischen gibt es mehr als 13.000 verschiedene	<input type="checkbox"/>	<input type="checkbox"/>



Best-Worst-Scaling NMT vs HT

	Best	Worst
Kryptowährungen. (Die meisten, seien Sie gewarnt, sind im Wesentlichen wertlos und werden es auch bleiben.)		
Die Einführung von Bitcoin löste eine Revolution in der Blockchain und Kryptowährung aus. Mittlerweile gibt es mehr als 13.000 verschiedene Kryptowährungen. (Die meisten, seien Sie gewarnt, sind im Wesentlichen wertlos und werden es auch bleiben.)	<input type="checkbox"/>	<input type="checkbox"/>

Q26

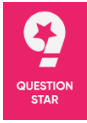
On a scale from **0 (lowest score)** to **4 (highest score)** , how good would you rate the translation you selected to be the **best** one?

- 4 (highest score)
- 3
- 2
- 1
- 0 (lowest score)

Q27

On a scale from **0 (highest score)** to **-4 (lowest score)** , how bad would you rate the translation you selected to be the **worst** one?

- 0 (highest score)
- 1
- 2
- 3
- 4 (lowest score)



Best-Worst-Scaling NMT vs HT

Q28

Please comment on what makes the one translation the best and the other the worst in your opinion. Feel free to add any other comments you would like to share on this text or the translations.

Q29

Set 4:

English original: You can think about energy efficiency like a budget. A financial budget tracks money as it's earned and spent. A budget for your body similarly tracks resources like water, salt, and glucose as you gain and lose them.

	Best	Worst
Sie können sich die Energieeffizienz wie ein Budget vorstellen. In einem Finanzbudget werden die Einnahmen und Ausgaben des Geldes erfasst. Ein Budget für Ihren Körper erfasst in ähnlicher Weise Ressourcen wie Wasser, Salz und Glukose, wie Sie sie zu- und abnehmen.	<input type="checkbox"/>	<input type="checkbox"/>
Am besten stellen Sie sich die Frage der Energieeffizienz so vor, als würden Sie ein Haushaltsbuch führen: Dabei notieren Sie, wie viel Geld hereinkommt und wie viel ausgegeben wird. Für Ihren Körper heißt das, dass Sie Ressourcen wie Wasser, Salz und Glukose eintragen und festhalten,	<input type="checkbox"/>	<input type="checkbox"/>



Best-Worst-Scaling NMT vs HT

	Best	Worst
wie viel Sie davon aufnehmen oder verbrauchen.		
Sie können sich Energieeffizienz wie ein Budget vorstellen. Ein Finanzbudget verfolgt, wie Geld verdient und ausgegeben wird. Ein Budget für Ihren Körper verfolgt in ähnlicher Weise Ressourcen wie Wasser, Salz und Glukose, während Sie sie gewinnen und verlieren.	<input type="checkbox"/>	<input type="checkbox"/>
Sie können sich Energieeffizienz wie ein Budget vorstellen. Ein Finanzhaushalt erfasst das verdiente und ausgegebene Geld. Ein Budget für Ihren Körper erfasst in ähnlicher Weise Ressourcen wie Wasser, Salz und Glukose, während Sie diese aufnehmen und verlieren.	<input type="checkbox"/>	<input type="checkbox"/>

Q30

On a scale from **0 (lowest score)** to **4 (highest score)** , how good would you rate the translation you selected to be the **best** one?

- 4 (highest score)
- 3
- 2
- 1
- 0 (lowest score)

Q31



Best-Worst-Scaling NMT vs HT

On a scale from **0 (highest score)** to **-4 (lowest score)**, how bad would you rate the translation you selected to be the **worst** one?

- 0 (highest score)
- 1
- 2
- 3
- 4 (lowest score)

Q32

Please comment on what makes the one translation the best and the other the worst in your opinion. Feel free to add any other comments you would like to share on this text or the translations.

Q33

Set 5:

English original: Physics in the late 19th and early 20th centuries was concerned with the spectrum of light being emitted by black bodies. A black body is a piece of material that radiates corresponding to its temperature — but it also absorbs and reflects light from its surroundings.

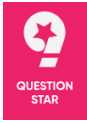
	Best	Worst
Die Physiker beschäftigten sich im späten 19. und frühen 20. Jahrhundert vor allem mit dem Lichtspektrum, das von sogenannten schwarzen Körpern ausgesendet wird. Ein schwarzer Körper ist ein Stoff, der wie alle anderen Körper seiner Temperatur entsprechend strahlt, aber auch Licht	<input type="checkbox"/>	<input type="checkbox"/>



Best-Worst-Scaling NMT vs HT

	Best	Worst
aus seiner Umgebung absorbiert und reflektiert.		
Die Physik des späten 19. und frühen 20. Jahrhunderts beschäftigte sich mit dem Lichtspektrum, das von schwarzen Körpern abgestrahlt wird. Ein schwarzer Körper ist ein Stück Material, das entsprechend seiner Temperatur strahlt - aber auch Licht aus seiner Umgebung absorbiert und reflektiert.	<input type="checkbox"/>	<input type="checkbox"/>
Die Physik des späten 19. und frühen 20. Jahrhunderts beschäftigte sich mit dem Spektrum des Lichts, das von Schwarzen Körpern emittiert wird. Ein schwarzer Körper ist ein Stück Material, das entsprechend seiner Temperatur strahlt – aber es absorbiert und reflektiert auch Licht aus seiner Umgebung.	<input type="checkbox"/>	<input type="checkbox"/>
Die Physik im späten 19. und frühen 20. Jahrhundert befasste sich mit dem Spektrum des von schwarzen Körpern emittierten Lichts. Ein schwarzer Körper ist ein Stück Material, das entsprechend seiner Temperatur strahlt – aber auch Licht aus seiner Umgebung absorbiert und reflektiert.	<input type="checkbox"/>	<input type="checkbox"/>

Q34



Best-Worst-Scaling NMT vs HT

On a scale from **0 (lowest score)** to **4 (highest score)** , how good would you rate the translation you selected to be the **best** one?

- 4 (highest score)
- 3
- 2
- 1
- 0 (lowest score)

Q35

On a scale from **0 (highest score)** to **-4 (lowest score)** , how bad would you rate the translation you selected to be the **worst** one?

- 0 (highest score)
- 1
- 2
- 3
- 4 (lowest score)

Q36

Please comment on what makes the one translation the best and the other the worst in your opinion. Feel free to add any other comments you would like to share on this text or the translations.

Q37

How easy/difficult was selecting the best translation in general for you?

- | | | | | |
|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|
| very easy | somewhat easy | neutral | somewhat difficult | very difficult |
| <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |



Best-Worst-Scaling NMT vs HT

Q38

How easy/difficult was rating the best translation for you?

very easy somewhat easy neutral somewhat difficult very difficult

Q39

How easy/difficult was selecting the worst translation in general for you?

very easy somewhat easy neutral somewhat difficult very difficult

Q40

How easy/difficult was rating the worst translation for you?

very easy somewhat easy neutral somewhat difficult very difficult

Q41

What do you think about this comparative approach of picking out the best and the worst translation?

Q42

What is your experience with assessing (human or machine) translation quality?

Appendix A. Questionnaire



Best-Worst-Scaling NMT vs HT

Q43

Was there anything about the survey you particularly liked or disliked? Please also add any other comments you would like to share on the topic.

Q44

You have reached the end of the survey. Thank you very much for taking the time to participate. Please click **End** in order to submit your responses.

Thank you!

Quantifying the Contribution of MWEs and Polysemy in Translation Errors for English–Igbo MT

Adaeze Ngozi Oluoba, Serge Sharoff, Callum Walker

Centre for Translation Studies,
School of Languages, Cultures and Societies
University of Leeds, LS2 9JT, UK

Abstract

In spite of recent successes in improving Machine Translation (MT) quality overall, MT engines require a large amount of resources, which leads to markedly lower quality for lesser-resourced languages. This study explores the case of translation from English into Igbo, a very low resource language spoken by about 45 million speakers. With the aim of improving MT quality in this scenario, we investigate methods for guided detection of critical/harmful MT errors, more specifically those caused by non-compositional multi-word expressions and polysemy. We have designed diagnostic tests for these cases and applied them to collections of medical texts from CDC, Cochrane, NCDC, NHS and WHO.

1 Introduction

In recent years, there has been increased research into improving the quality of machine translation (MT) outputs (Wu et al., 2016; Hassan et al., 2018). Evidenced by the switch from rule-based and statistical MT systems to neural MT systems, this has led to visible improvement of MT outputs. However, these improvements are more common with ‘high-resourced languages’ that have sufficient data resources for training MT models. Thus, ‘low/under-resourced languages’ like Igbo, lag behind in this progress. The Igbo language (Ásusu Ìgbò) is one of the three major languages spoken in Nigeria, it is the native language of the Igbo people, an ethnic group in South-Eastern Nigeria. It is

also, a recognised minority language in Equatorial Guinea and Cameroon¹.

Regarding language resources, Igbo language, for instance, has only 18,369 Wikipedia articles as of 02 October 2023 unlike the English and French languages that are in the top 5 languages used in Wikipedia with 6,722,185 and 2,557,357 articles respectively. Additionally, there are no single parallel corpora with Igbo language as a language pair in Sketch Engine² and only a few available in OPUS³. The amount of parallel data that includes Igbo language as one of the language pairs remains limited. Thus, “Igbo – any language” is a low-resource language pair.

Another critical aspect of the research into MT output is the evaluation of MT for health domains. This is especially important given the recent experience with the COVID-19 pandemic where the majority of the world’s population had to be confined in isolation. Prior to the COVID-19 pandemic, there had been calls for increased research into translation in time of crisis championed by the International Network on Crisis Translation.⁴ One of the goals of this network was “to make meaningful and effective contributions ... that enable accurate and timely translation-enabled crisis communication, with a particular focus on health-related content”.

In crisis, access to information in one’s L1 cannot be over-emphasized, and given the speed, cost-effectiveness, and easy availability of MT during crisis or in situations like the self-isolation necessitated by the coronavirus pandemic, it is safe to

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹<https://www.africanexponent.com/8-most-spoken-local-languages-in-africa/>

²<https://app.sketchengine.eu/>

³<https://opus.nlpl.eu/>

⁴<https://cordis.europa.eu/project/id/734211>

assume as pointed out by O'Brien (2022), that MT would be the most logical 'go-to' tool for users. This has been evidenced many times, the most recent being the massive deployment of MT during the Ukraine crisis (Cirule, 2022)

Regarding the reliability of these MT tools, there have been recent claims comparing MT quality to be relatively close to human translations (Wu et al., 2016; Hassan et al., 2018). However, there are still some reservations on the perceived high quality of machine translation based on some qualitative evaluations (Läubli et al., 2020; Wang et al., 2021). These evaluations, which were carried out on news texts, found that translating between languages with varying word orders posed a challenge for machine translation systems. They also reported MT systems' weak error tolerance, which makes them susceptible to inaccurate translations due to minor punctuation or spelling errors that might be overlooked by a human. The inability to discern what should be emphasised or omitted; and data sparsity in low-resource language pairs and domains were also identified as hindrances to human-machine parity in machine translation. Dew et al. (2018) noted that for statistical-based machine translation systems, language resources matter, as they are known to perform better with language pairs that are well represented online. Other studies on neural machine translation systems (Lakew et al., 2018; Murthy et al., 2018; Singh and Singh, 2022) also agree that MT quality into low-resourced languages is relatively low.

In his empirical evaluation of the quality of machine translation of 20 phrases from English into 107 languages, using Google Translate (GT),⁵ Benjamin (2019b)'s study recorded good, almost perfect outputs for high-resourced languages and lesser quality outputs for low-resourced languages. For Igbo language, he reported that 47.5 per cent of the texts were accurately translated while noting that GT was able to provide fairly meaningful translations 60 per cent of the time. Even so, his study was on short non-ambiguous phrases of common usage.

More studies to evaluate and improve MT quality for health domain and especially into/from low-resourced languages is therefore considered imperative as the accuracy of health information received in times of crisis is vital. Our aim here is to provide evaluation of specific phenomena for auto-

matic identification of translation problems. As such, an integral part of this research includes highlighting problem areas that negatively affect MT quality of health-related texts from English – Igbo, as this area has not yet been explored.

In our preliminary studies, we have identified MWEs and polysemy as the most common causes of critical errors.

Research questions:

- What proportion of critical errors can be identified via automatic detection of MWEs and polysemy?
- What is the most appropriate granularity for error detection: sentence, segment or full text?

2 Related Work

Multi-Word Expressions (MWEs): Multi-word expressions (MWEs) have been defined by Sag et al (2002) as a combination of words for which syntactic or semantic properties of the whole expression cannot be obtained from its parts. They are lexicalized combinations of two or more words that are exceptional enough to be considered as single units in the lexicon (Schneider et al., 2014). Non-compositionality in phraseology (difficulty in deriving meaning from individual components) and non-substitutability (components cannot be replaced with synonyms) have been reported by Premasiri and Ranasinghe (2022) as features of MWEs that are challenging for NLP (Natural Language Processing).

Zaninello and Birch (2020) in evaluating the effect of annotation and data augmentation in the English – Italian translation of MWEs in a neural machine translation system, note that for non-compositional MWEs, the translation quality was especially low. They report that following their study, there is clear indication that NMT systems find it difficult to translate non-compositional MWEs even for high-resourced languages, and that focusing on improving MWEs in a text can not only improve the quality of translation of MWEs in the text but also the overall machine translation quality. Arvi (2018) in his comparison of a rule-based machine translation system, (SALAMA), and Google Translate's translations of multi-word expressions in news texts from English – Swahili, discovered that Salama performed better than GT in translating MWEs. He therefore advocated the

⁵<https://translate.google.com/>

investment into rule-based system for translating highly inflectional low-resourced languages, as the rules can be adapted to similar languages and the accuracy of the translation would not be dependent on large parallel data.

Polysemy: Abdelaal and Alazzawie (2020) posit that Arnold et al. (1994)'s stance that ambiguity is a big challenge for MT no longer holds true for Google Translate since its switch to a neural system. Nevertheless, Xie et al (2021) report the occurrence of inconspicuous yet clinically significant medical and health (English–Chinese) machine translation mistakes suspected to be due to the limited ability of current neural MT systems to correctly interpret the meaning of polysemous words. Thus, leading to an increase in risks for end-users of machine translation systems. Meenal and Govindarajan (2023) whose research was on the challenge of machine translating polysemous words across various domains from French–English on Google Translate, concluded that the translations confirmed MT's current incapability to correctly translate polysemy even for a high-resourced language pair like French and English.

The above cases show that polysemy is a challenge for MT in high-resourced languages. This is also the case for low-resourced languages as seen in other studies. For instance, in their evaluation of machine translated texts into Lithuanian across two MT systems, Google Translate and VDU, Petkevičiūtė and Tamulynas (2011) report that the two systems recorded similar significant challenges in their translation of polysemous words. Likewise, Tudor (2017)'s research on machine translating polysemous Croatian words into English language, revealed a low level of translation accuracy. Abdulaal (2022) also reports that polysemous words should be considered while machine translating literary texts from English to Arabic as the texts could contain errors due to the machine translation system's inability to properly translate such linguistic phenomenon.

3 Methodology

3.1 Detection of MWEs

There have been a few studies on the detection and classification of MWEs for use in NLP. Zaninello and Birch (2020) report that they used manually compiled entries from a bilingual and a monolingual dictionary, instead of an automatic tool to ex-

tract MWEs in order to maximise accuracy during the extraction process. Simkó et al (2017) used POS tagging and dependency parsing to detect verbal MWEs.

For MWE identification, we treated them on the basis of syntax. We noted that our data (see 3.3) contained terminological units, named entities and light verb constructions. Thus to detect and extract MWEs in our study, Spacy's POS tagger and dependency parser were used to identify syntactic patterns on the texts. Due to the multi-word named entities in the texts, we did not apply n-gram restrictions during the extraction process. After tagging the texts, we subsequently extracted the error bearing segments and manually tagged them as a test set to compare Spacy's accuracy. We determined a precision and recall score of 0.91 each, which corresponds to Spacy's accuracy evaluation claim.

3.2 Detection of Polysemy

Analysis of the Corpus of Contemporary American English shows that most common English words have at least two senses, which produces 50/50 odds in the possible case that the target language uses different words for those different senses (Benjamin, 2019a). The word "back" for instance is reported to have 36 different senses, multiplying its translation possibility by 36 if each sense correlates to a different word in a target language. There are at least 6 different translations for "back" in Igbo ("azu", "nkezu" "nke gara aga", "n'azu", "ikwado", "ebe azu"). Amongst these six translations, "azu" can also mean at least 8 different things; (back, fish, train, behind, bum, shark, retreat, rear). Scenarios like this could lead to what Benjamin (2019a) describes as the multiplication effect of polysemy in translation.

To identify polysemous words in our corpus, we used NLTK's WordNet Interface (Miller, 1990), to identify the number of their word senses. We also incorporated the use of domain statistics (Hamilton et al., 2016) and used Word2vec-google-news-300 to identify the number of contexts/domains a polysemous word can occur in. This we did to determine if the number of word senses and number of contexts a word has, affects English–Igbo MT accuracy.

3.3 Data

We collated a total of 123 English texts, approximately 200,000 words, from the United States Cen-

ters for Disease Control and Prevention (CDC)⁶, Cochrane⁷, the Nigeria Centre for Disease Control (NCDC)⁸, the United Kingdom's National Health Service (NHS)⁹ and the World Health Organization (WHO)¹⁰ websites. These texts, published between 2019 and 2022, are primarily about COVID-19 and are either instructional or informative texts with the exception of the Cochrane text which is majorly professional and academic. The first phase of this research involved a preliminary study. For this preliminary study, we selected one text from each source, comprising a total of 168 sentences and 2000 words. We thereafter grouped the selected data into two, considering variety in terminology. Flesch reading ease score was used to assess the linguistic difficulty of the English texts:

a) Reference Information texts for Health Professionals (henceforth 'PROF'): This text contained a lot of medical terminology. The Flesch-Kincaid score for this text is 27.2 and the Flesch-Kincaid Grade level is 13.2, and thus classed as very difficult to read for non-professionals. 28 of our selected sentences had this classification.

b) Instructional and Informative texts for Public/Patients (henceforth 'Info'): The Flesch-Kincaid reading ease score for this text is 57.2 and the Flesch-Kincaid Grade level is 8.5. The text is classed as a simple text with simple syntax. It is intended to be informative for the public and there is minimal use of highly specialised terminology and acronyms. Some part of the text also contains instructions. There were 140 sentences in the "Info" text.

Our study is based on output from Google Translate as in our preliminary evaluation of three MT systems, it emerged as the best tool for the English-Igbo language pair.

Given our aim to find which detection parameters from Sections 3.1 and 3.2 provides the most appropriate granularity for error-detection, we also vary the window for detection from tokens to segments and to full texts. For the word-level evaluation, our data contained 1490 tokens for the Info text and 388 for the PROF text. We also divided the texts into segments of meaning; 171 segments for the Info text and 45 segments for the PROF text.

⁶<https://www.cdc.gov/>

⁷<https://www.cochrane.org/>

⁸<https://ncdc.gov.ng/>

⁹<https://www.nhs.uk/>

¹⁰<https://www.who.int/>

3.4 Annotation Guidelines and Classification of Error Categories

Annotation guidelines were prepared to improve uniformity of MT quality evaluation across the texts. Given the absence of parallel data for our selected texts, we human translated the selected English texts into Igbo, as a gold standard for evaluation. For the error classification, we combined both linguistic and medical errors; linguistic errors here are errors that border on language fluency whereas medical errors refer to errors that though being linguistically fluent, contain errors that are medically significant and can cause harm. We thus grouped the MT errors into three error categories vis general errors, syntactic errors and terminology errors. So, if an error is neither a terminology nor syntactic error, it is tagged as a general error. Thereafter, if there are a lot of major and critical errors that have been classed as general errors, the error-causing words will then be further analysed. This phase aims to confirm Xie et al (2021)'s claim that terminology is not the major challenge in machine translation of health texts and also confirm if major and critical syntactic errors are made during English-Igbo MT. For the three error classes, if a segment had more than one error category, the category with the higher error severity was applied.

An error penalty was also associated to each error category according to the severity of the error. We carried out a three-level assessment scale for this study by modelling the error severity guidelines described by O'Brien (2012) and Comparin and Mendes (2017) which were adopted from the Multidimensional Quality Metrics (MQM) framework and Localization Quality Evaluation (LQE). Following the description of each level (written below), we found the three-level assessment to be a good fit for the preliminary error analysis. Error categories/potential for harm from inaccurate machine translations were thus grouped into three levels (Minor/no potential for harm, Major/potential for harm, Critical errors/life-threatening/ catastrophic/harmful). We also favoured an arithmetic progression of error penalty score (1,2,and 3) in place of the geometric scores used by Comparin and Mendes (2017) because we wanted to make linear distinctions among the error categories, thus making the difference between categories easier to interpret and apply consistently.

i) Minor: linguistically, the output is wrong, but the reader can decode the meaning of the sentence;

medically, the output is wrong but does not affect understanding nor cause any harm. Minor errors have a score of 1.

ii) Major: linguistically, this is a wrong output that hinders the understanding of the text; medically, the output can cause a degree of harm that is not life-threatening. Major errors carry a weighting of two points.

iii) Critical errors: errors that make the text incomprehensible, can cause harm or that connote meaning that is opposite of the source text. 3 points are assigned to each critical error. For instance:

- Source Text (ST): even if you've had a **positive** test result for COVID-19 before.

MT: ọbụlagodi na ị nwetala nsonaazụ nyocha **dị mma** maka COVID-19 na mbụ.

Back Translation (BT): Even if you have had a **negative** test result for COVID-19 before.

This segment was annotated as a terminological error and had an error score of "3" because the error is a critical error and could cause life threatening harm.

- ST: If you **book for** someone:...

MT: Ọ buru na ị **na-edede akwụkwọ maka** mmadu:

BT: If you are **writing about** someone

This segment containing a linguistic error was annotated as a general error with an error score of "2".

- ST: Have a high **temperature**

MT: Nwee **okpomọkụ** dị elu

BT: Have a high **hotness**

This segment was annotated as a terminological error with an error score of "1" as the meaning can be deduced.

Quantifying Error Severity

A cross-genre similarity is also identified, as error severity follows the same pattern for both the PROF and Info texts. The results in Table 1 indicate that there are more major and critical general errors, causing about 51 per cent of total errors annotated. Terminology based major and critical errors constitute about 24 per cent of total errors. We find that the MT system did not record any major or critical syntactic errors, as the only

syntactic errors were minor and did not distort the meaning of the text. This conforms with Xie et al. (2021)'s opinion, that terminology is not the major challenge in machine translation of health texts which we sought to test its applicability to English-Igbo machine translation and thus decide if a fine-grained analysis on the exact typology of these errors would be beneficial. Subsequent to this finding, the next section is dedicated to an in-depth analysis of these errors.

3.5 Analysis of Major and Critical MT Errors

In this phase, we merge major and critical errors into one label (harmful) thus narrowing the criteria to either a negligible error or a harmful error. We also run further in-depth harmful-error analysis on the level of tokens. We therefore used Spacy's POS tagger to identify the parts of speech and dependencies of the tokens that caused harmful MT errors.

As evident in Table 2, adverbs and adjectives cause the highest amount of harmful errors in the Info text, whereas nouns and verbs account for the most errors in the PROF text (Table 3). Given the fact that the PROF text contains a lot of medical terminology which are mainly nouns and verbs and the Info text has a lot of descriptions as an informative text, we find that these distinctive features contribute to the difference in the ranking of the top 5 error-causing parts of speech for the two texts. However, a notable similarity in this scenario is that the error causing parts of speech for the Info also features in the top 5 harmful error causing parts of speech for the PROF text [ADV, ADJ, NOUN, VERB, PROP]. We therefore analyse the features of the error-causing words to determine if there are more definite similar features between the cause of harmful errors for both text types.

3.6 Multi-word Expressions and their Impact as Cause of Critical Errors in English-Igbo MT

One other distinctive feature of the error-causing POSs was that they formed part of multi-word expressions. This part of the experiment thus served to reveal if the top error-causing parts of speech for both text types have similar linguistic classifications as part of MWEs.

Thus, to analyse the effect of multi-word expressions on English-Igbo MT, we would focus on their frequency, syntactic constructions and semantic properties in the source text.

Text Type	General Errors			Syntactic Errors			Terminology Errors		
	Negligible**	Harmful***		Negligible	Harmful		Negligible	Harmful	
	Min.	Maj.	Crt. *	Min.	Maj.	Crt.	Min.	Maj.	Crt.
PROF	10	13	8	1	0	0	0	10	5
Info	24	38	26	5	0	0	2	13	12

Table 1: Error severity by segments. *Critical errors **Negligible: scores < 2 *** Harmful errors: scores ≥ 2

POS	Freq.	Harmful Error	% Error
ADV	71	15	21.13%
ADJ	118	24	20.34%
NOUN	330	45	13.64%
VERB	246	29	11.79%
PROPN	72	8	11.11%
AUX	72	5	6.94%
SCONJ	50	3	6.00%
ADP	161	9	5.59%
DET	78	4	5.13%
PRON	149	7	4.70%
PART	52	2	3.85%
CCONJ	73	2	2.74%
NUM	14	0	0.00%
X	2	0	0.00%

Table 2: Frequency of Harmful-Error causing POS (Info)

POS	Freq.	Harmful Error	% Error
NOUN	127	40	31.50%
VERB	47	13	27.66%
PROPN	20	4	20.00%
ADJ	34	6	17.65%
ADV	10	1	10.00%
ADP	43	3	6.98%
CCONJ	26	0	0.00%
DET	22	0	0.00%
AUX	16	0	0.00%
PRON	15	0	0.00%
PART	10	0	0.00%
SCONJ	8	0	0.00%
NUM	6	0	0.00%

Table 3: Frequency of Harmful-Error causing POS (PROF)

Frequency Distribution of MWEs in Source Text:

Further analysis of the syntactic information of the error causing POSs and their collocates highlights that they form part of MWEs. Table 6 shows the frequency of these MWEs and the percentage errors caused by the MWEs for both text types.

Syntactic Properties of Harmful Error-causing MWEs:

Classes of error-causing MWEs were distinguished based on their categorical properties and their syntactic features. It is note-worthy that most compounds in the data used are open (i.e. the compound words are written with spaces), non-

compositional and non-hyphenated. This could contribute to a machine translation system’s inability to properly distinguish its linguistic features and give an accurate translation. This was the case for 38 per cent of open compounds in the PROF text.

Using Spacy’s syntactic dependency parser, we also analysed the dependency tags (Schuster and Manning, 2016) of the words that caused harmful errors in the MT output.

Dep	Count	Error	% Error
quantmod	5	2	40%
oprd	3	1	33%
acomp	20	5	25%
appos	4	1	25%
amod	99	20	20%
ccomp	28	5	18%
advmod	88	15	17%
nsubjpass	6	1	17%
doj	141	23	16%
xcomp	31	5	16%
relcl	19	3	16%
npadvmod	13	2	15%
pcomp	14	2	14%
pobj	131	17	13%
advcl	53	6	11%
auxpass	10	1	10%
conj	114	10	9%
compound	81	7	9%
neg	12	1	8%
acl	13	1	8%
mark	29	2	7%
prep	152	9	6%
det	77	4	5%
aux	80	4	5%
nsubj	92	4	4%
cc	73	2	3%

Table 4: Error Frequency > 1 by Dep. (Info)

We see in Table 4 and Table 5 that modifiers e.g. “quantmod”, “amod”, and complements e.g. “acomp”, “xcomp” which form part of noun phrases, compound nouns and verb phrases are frequent causes of harmful errors.

Noun Phrases (NP) and Compound Nouns:

Syntactically, compound nouns in English are typically left-branching i.e., the modifiers come before the noun whereas the reverse is the Igbo language case (right-branching: the modifiers come after the

Dep	Count	Error	% Error
prt	1	1	100%
acl	4	3	75%
npadvmod	3	2	67%
xcomp	4	2	50%
conj	23	10	43%
compound	34	14	41%
amod	31	11	35%
pobj	41	12	29%
dobj	25	5	20%
relcl	6	1	17%
advmod	13	2	15%
advcl	8	1	13%
prep	39	2	5%
ROOT	21	1	5%

Table 5: Error Frequency >1 by Dep. (PROF)

noun) e.g., hand-cream: ude aka (ude=cream, aka=hand). For noun phrases, English language accepts both forms of modification whereas Igbo accepts only post-modification. (Orji et al., 2022).

We find that 33 and 44 per cent of compound nouns and noun phrases in the Info and PROF text are causes of harmful errors, these errors nevertheless were not caused by the syntactic difference between compound nouns and noun phrases in English and Igbo (see Table 6).

Compound Verbs, Verb Phrases (VP) and Light Verb Constructions (LVC): Syntactically, compound verbs in English are typically right-branching i.e. the modifiers come after the verb, however English language accepts both forms of modification e.g. double-click (pre-modified) or throw up (post-modified). Whereas in Igbo, compound verbs are strictly post-modified e.g. “weta = we-ta [to bring]”. Despite this syntactic difference between English and Igbo, we did not record any harmful compound verb/VP/LVC syntactic translation error in both texts. Babych et al. (2009)’s observation that rule-based MT often mistranslates LVCs, still holds true in this English–Igbo neural MT experiment as seen in Table 7 and this mistranslation was caused by the semantic implication of LVCs.

Multi-Word Named Entities: Our experiments proved single word named entities did not cause any harmful errors. Multi-word named entities, on the other hand, were responsible for some harmful errors (20 per cent in PROF text). Thus, the need to have (multi-word) named-entities as part of the multi-word expressions. Results from our syntactic analysis of multi-word named entities further proved that the major challenge of English–Igbo

MT is not primarily syntactic as Google Translate did not output any significant syntactic errors in its translation of MWEs.

Semantic Properties of Harmful-Error Causing MWEs: Given the results above, we thereafter investigated the semantic properties of the MWEs in our data. Dickins (2020) defined multi-word expressions in relation to their semantic compositions. He also classified them into three viz: “Type 1: fully non-compositional, i.e. none of the words has an independent sense; Type 2: at least one of the words has a sense which is independent but is only found in the context of this expression; and Type 3: at least one of the words has a sense which is independent but is only found in definable limited contexts of which this context is one.”

A greater percentage of the MWEs in our data are open compounds and endocentric or copulative, this corresponds with Dickins (2020)’s type 2 and 3 MWEs. Less than 5 per cent of our dataset contained closed compounds (i.e. the compound words are written with no spaces or punctuation) and these closed compounds did not cause any harmful errors. One other semantic feature of note is that some harmful-error causing multi-word expressions which are type 2 and type 3 compounds (independent contextual senses) contained individual polysemous words e.g. ‘positive test result’. This will be discussed in the next section.

3.7 Polysemy and its Impact as a Cause of Critical Errors in English–Igbo MT

Collocational relations and context are meant to be helpful in neural machine translation systems; nonetheless, polysemy is one of the linguistic phenomena that has been noted as a challenge to MT especially when the probability of the accurate translation of the word in context is statistically low i.e. not the most frequent sense, or its sense is insignificant in the MT system’s training data for the languages in contact.

Error severity by word senses: In Table 8 and Table 9, we record the frequency and severity of the polysemous words in the data by their word-senses. We investigated if the error-causing rate of a polysemous word is directly proportional to the number of word-senses it has. One constant is that there is a similar frequency in the percentage of the errors/harmful errors caused by polysemous words in both text types. Furthermore, in at least 76 per cent of the time across all word-senses and text

MWE	Freq Info	Errors	Error %	Freq Prof	Errors	Error %
Compound Nouns/NP	82	27	33%	41	18	44%
Compound Verbs	23	9	39%	1	1	100%
Multi-Word Named Entities	24	1	4%	5	1	20%

Table 6: Manually annotated MWE errors in both texts

Source Text	Machine Translation	Back Translation
Take a break	Were ezumike	Collect a break
Get Vaccinated	Were ogwu mbochi	Collect a vaccine

Table 7: Example cases of inaccurate machine translation of VPs and LVCs

Info	≥ 2	≥ 5	≥ 7	≥ 10	≥ 15	≥ 20
Frequency	807	418	307	227	133	67
Error (201)	133	98	70	50	30	13
Harmful Error (153)	106	82	57	38	24	10
% of error is polysemous	66%	49%	35%	25%	15%	6%
% of harmful error is polysemous	69%	54%	37%	25%	16%	7%
% of Polysemous word is a harmful error	13%	20%	19%	17%	18%	15%
% of polysemous error is harmful	80%	84%	81%	76%	80%	77%

Table 8: Error severity by word senses (Info)

Prof	≥ 2	≥ 5	≥ 7	≥ 10	≥ 15	≥ 20
Frequency	220	105	80	43	24	3
Error (82)	61	40	32	14	5	1
Harmful Error (67)	55	39	31	14	5	1
% of error is polysemous	74%	49%	39%	17%	6%	1%
% of harmful error is polysemous	82%	58%	46%	21%	7%	1%
% of Polysemous word is a harmful error	25%	37%	39%	33%	21%	33%
% of polysemous error is harmful	90%	98%	97%	100%	100%	100%

Table 9: Error severity by word senses (PROF)

types, the polysemous error is a harmful error. This reveals that polysemous words do not just cause MT errors; they cause harmful errors in English – Igbo machine translation of medical texts. Another important point is that the percentage of errors for polysemous words of word-senses greater than 10 is comparatively lower than words of word-senses 7 and below.

Error severity by Polysemy domain/context: This part of the study sought to analyse if the error-causing rate of a polysemous word is directly proportional to the number of domains or contexts (con) it occurs in. We thus varied our experiments to account for different context lengths; words occurring in greater than “1,2,5, and 10” domains/context. However, the results show that at least 30 per cent of the polysemous words in both

the Info and PROF texts were causes of harmful errors irrespective of the number of contexts the polysemous word has. For the PROF text, all the polysemous words that had up to ten contexts caused not just errors but harmful errors (Table 10 and Table 11).

Below are examples of errors caused by polysemous words in this study.

i) Source Text (ST): Always call before **visiting** your doctor or health facility.

Machine Translation (MT): Na-akpọ oku mgbe niile tupu **iga leta** dokita gi ma o bu ulo oru ahike.

Back Translation (BT): Always call before you **pay a social visit** to your doctor or health facility.

The idea here is one of going to a health centre to be seen by the health professional, not a ‘social visit’ as translated.

No of Domain	≥1	≥2	≥5	≥10
Frequency	213	167	46	10
Error (201)	97	78	20	4
Harmful Error (153)	75	63	15	3
% of error is con	48%	39%	10%	2%
% of harmful error is con	49%	41%	10%	2%
% of con word is harmful error	35%	38%	33%	30%
% of con error is harmful	77%	81%	75%	75%

Table 10: Error severity by no. of domains (Info)

No of Domain	≥1 (77)	≥2 (56)	≥5 (17)	≥10 (5)
Frequency	77	56	17	5
Error (82)	48	34	10	5
Harmful Error (67)	42	29	9	5
% of error is con	59%	41%	12%	6%
% of harmful error is con	63%	43%	13%	7%
% of con word is a harmful error	55%	52%	53%	100%
% of con error is harmful	88%	85%	90%	100%

Table 11: Error severity by no. of domains (PROF)

ii) ST: ...such as the emergency department or **dedicated** COVID-19 clinics.

MT: *dị ka ngalaba mberede ma ọ bụ ụlọ ọgwụ COVID-19 raara onwe ya nye.*

BT: Such as the emergency unit or COVID-19 hospital that has **committed itself**.

The translation of 'dedicated' here is that of a person instead of an allocated item.

iii) ST: ... at both title and abstract, and full-text stage.

MT: ... ma aha ma nke nkịti na **okwa** ederere zuru oke.

BT: ... at both title and normal, and full-text **podium**. *'Stage' is translated as a theatre stage instead of its accurate connotation of a process.*

iv) ST: **Cough** or **sneeze** into a **tissue**.

MT: **Ukwara** ma ọ bụ **uzere** n'ime **anụ ahụ**.

BT: **A cough** or **a sneeze** into the **body**.

'Tissue' translated as 'anụ ahụ': body tissue, instead of its implied context of 'tissue paper'

v) ST: Talk about your concerns – anxiety at this time is **normal**.

MT: Kwuo banyere nchegbu gị - nchegbu n'oge a bụ **ihe nkịti**.

BT: Talk about your concerns- anxiety at this time is **insignificant**.

Here the polarity of 'normal' is misinterpreted. The translation does not adequately represent the sentiments expressed. It unfortunately stifles the

emotions of anxiety.

vi) ST: even if you've had a **positive** test result for COVID-19 before.

MT: ọbụlagodi na ị nwetala nsonaazụ nyocha **dị mma** maka COVID-19 na mbụ.

BT: Even if you have had a **negative** test result for COVID-19 before.

4 Discussion

Ambiguity in MWEs: Even though most of the MWEs in the source texts were endocentric, 50 per cent of the MWEs in the PROF text contained at least one polysemous word (type 3 MWE) which posed a challenge for MT and made the semantics not easily predictable from the expression. This caused an error in at least 64 per cent of MWEs with polysemous constituent words. Examples ii, iii and vi above highlight some of the cases. Google Translate was unable to recognise cases in which expressions with seemingly positive connotations are used for expressing a negative idea e.g., positive covid-19 test result in example (vi) was translated as 'nyocha dị mma maka COVID-19' implying a negative test result, as the MT system uses the connotation that 'positive' implies something good. This is different to its medical meaning showing the presence of an organism/disease. This reveals that polysemous words and multi-word expressions are to be analysed independently as pol-

polysemous words can be part of MWEs but not vice versa. The results from our study can also aid in evaluation-guided pre-editing (Babych et al., 2009) for English – Igbo machine translation and the resulting MT output could be re-evaluated to quantify pre-editing impact.

5 Conclusion and Future Work

In our paper, we have identified and quantified what linguistic features of English as a source language, create challenges for a machine translation system to accurately translate a medical text into Igbo. Our findings confirm that a medical text filled with multi-word expressions and polysemous words is not suitable for English to Igbo machine translation as Google Translate is still unable to correctly translate such linguistic properties from English to a low-resource language like Igbo. We also find that syntactic differences between the two languages do not contribute to harmful MT errors. For polysemous words, focusing on their word senses reveals an error- peak point of seven word senses, whereas all levels of domain/context numbers had a high percentage of harmful errors. As such, future work to determine if these challenges will still persist on a larger data set will be primarily on word-senses less than and equal to seven and there will be no focus on the number of contexts. Token- level analysis of our data resulted in more detailed findings and would also form the guideline of further work. As part of our wider objectives, we intend to use the results from this preliminary study to develop machine learning classifiers in order to predict medical texts that could be catastrophic for MT users to machine translate from English to Igbo. Finally, we hope our findings can also serve as a guide to evaluate/detect causes of critical MT errors for low-resourced languages especially other Niger-Congo language families.

Acknowledgements: We thank the anonymous reviewers for their time and their insightful suggestions and comments.

References

Abdulaal, Nouredin Mohamed and Abdulkhaliq Alazawie. 2020. Machine translation: The case of Arabic-English translation of news text. *Theory and Practice in Language Studies*, Vol. 10(No. 4):408–418.

Abdulaal, Mohammad Awad Al-Dawoody. 2022. Tracing machine and human translation errors in some

literary texts with some implications for EFL translators. *Journal of Language and Linguistic Studies*, 18(Special Issue 1):176–191.

Arnold, Doug, Lorna Balkan, Siety Meijer, R. Lee Humphreys, and Louisa Sadler. 1994. *Machine Translation: an Introductory Guide*.

Arvi, Hurskainen, 2018. *Sustainable language technology for African languages*, book section Sustainable language technology for African languages. Routledge.

Babych, Bogdan, Anthony Hartley, and Serge Sharoff. 2009. Evaluation-guided pre-editing of source text: improving MT-tractability of light verb constructions. In *Proc 13th European Association for Machine Translation*, pages 36–43, Barcelona, May.

Benjamin, Martin. 2019a. The astounding mathematics of machine translation, 01 April 2019.

Benjamin, Martin. 2019b. Empirical evaluation of Google Translate across 107 languages, 2019.

Cirule, Gunta. 2022. Tilde has developed Ukrainian machine translation systems to help refugees.

Comparin, Lucia and Sara Mendes. 2017. Using error annotation to evaluate machine translation and human post-editing in a business environment.

Dew, Kristin N., Anne M. Turner, Yong K. Choi, Alyssa Bosold, and Katrin Kirchhoff. 2018. Development of machine translation technology for assisting health communication: A systematic review. *Journal of Biomedical Informatics*, 85:56–67.

Dickins, J. 2020. An ontology for collocations, formulaic sequences, multiword expressions, compounds, phrasal verbs, idioms and proverbs. *Linguistica Online*, 23:29–72.

Hamilton, William L., Jure Leskovec, and Dan Jurafsky. 2016. Diachronic word embeddings reveal statistical laws of semantic change. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1489–1501, Berlin, Germany, August. Association for Computational Linguistics.

Hassan, Hany, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, Shujie Liu, Tie-Yan Liu, Renqian Luo, Arul Menezes, Tao Qin, Frank Seide, Xu Tan, Fei Tian, Lijun Wu, and Ming Zhou. 2018. Achieving human parity on automatic Chinese to English news translation.

Lakew, Surafel M., Marcello Federico, Matteo Negri, and Marco Turchi. 2018. Multilingual neural machine translation for low-resource languages. *Italian Journal of Computational Linguistics*, 4(1):11–25.

- Läubli, Samuel, Sheila Castilho, Graham Neubig, Rico Sennrich, Qinlan Shen, and Antonio Toral. 2020. A set of recommendations for assessing human-machine parity in language translation. *The Journal of Artificial Intelligence Research*, 67:653–672.
- Meenal, T S and P Govindarajan. 2023. The challenges of using machine translation while translating polysemous words. 2023, 11(2):5.
- Miller, George. 1990. WordNet: an online lexical database. *International Journal of Lexicography*, 3(4).
- Murthy, Rudra, Anoop Kunchukuttan, and Pushpak Bhattacharyya. 2018. Addressing word-order divergence in multilingual neural machine translation for extremely low resource languages.
- O'Brien, Sharon. 2012. Towards a dynamic quality evaluation model for translation. *The Journal of Specialised Translation*, (17):55–77.
- O'Brien, Shanon. 2022. Crisis translation: A snapshot in time. *INContext: Studies in Translation and Interculturalism*, 2(1).
- Orji, Ifeoma Maryann, Sylvanus Okwudili Anigbogu, Oluchukwu Uzoamaka Ekwealor, and Ukamaka Bertrand Chidi. 2022. Enhanced machine learning algorithm for translation of English to Igbo language. *Machine Learning Research*, 7(1):8–14.
- Petkevičiūtė, Inga and Bronius Tamulynas. 2011. Computer-based translation into lithuanian: Alternatives and their linguistic evaluation. *Studies about Languages*, (18):38–45.
- Premasiri, Damith and Tharindu Ranasinghe. 2022. BERT(s) to detect multiword expressions. *ArXiv*, abs/2208.07832.
- Sag, Ivan, Timothy Baldwin, Francis Bond, Ann Copestake, and Dan Flickinger. 2002. *Multiword Expressions: A Pain in the Neck for NLP*.
- Schneider, Nathan, Spencer Onuffer, Nora Kazour, Emily Danchik, Michael T. Mordowanec, Henrietta Conrad, and Noah A. Smith. 2014. Comprehensive annotation of multiword expressions in a social web corpus. Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14), pages 455–461. European Language Resources Association (ELRA).
- Schuster, Sebastian and Christopher D. Manning. 2016. Enhanced English Universal Dependencies: An improved representation for natural language understanding tasks. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 2371–2378, Portorož, Slovenia, May. European Language Resources Association (ELRA).
- Simkó, Katalin Ilona, Viktória Kovács, and Veronika Vincze. 2017. USzeged: Identifying verbal multiword expressions with POS tagging and parsing techniques. In *Proceedings of the 13th Workshop on Multiword Expressions (MWE 2017)*, pages 48–53, Valencia, Spain, April. Association for Computational Linguistics.
- Singh, Salam Michael and Thoudam Doren Singh. 2022. Low resource machine translation of English–Manipuri: A semi-supervised approach. *Expert systems with applications*, 209.
- Tudor, Atena. 2017. *Machine Translations of Polysemous Croatian Words in Various Text Genres*. Thesis.
- Wang, Haifeng, Hua Wu, Zhongjun He, Liang Huang, and Kenneth Ward Church. 2021. Progress in machine translation. *Engineering*.
- Wu, Yonghui, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation.
- Xie, Wenxiu, Meng Ji, Riliu Huang, Tianyong Hao, and Chi-Yin Chow. 2021. Predicting risks of machine translations of public health resources by developing interpretable machine learning classifiers. *International Journal of Environmental Research and Public Health*, 18(16):8789.
- Zaninello, Andrea and Alexandra Birch. 2020. Multiword expression aware neural machine translation. Proceedings of the Twelfth Language Resources and Evaluation Conference, pages 3816–3825. European Language Resources Association.

Analysis of the Annotations from a Crowd MT Evaluation Initiative: Case Study for the Spanish-Basque Pair

Nora Aranberri

HiTZ Basque Center for Language Technologies - Ixa NLP Group
University of the Basque Country (UPV/EHU)
nora.aranberri@ehu.eus

Abstract

With the advent and success of trainable automatic evaluation metrics, creating annotated machine translation evaluation data sets is increasingly relevant. However, for low-resource languages, gathering such data can be challenging and further insights into evaluation design for opportunistic scenarios are necessary. In this work we explore an evaluation initiative that targets the Spanish—Basque language pair to study the impact of design decisions and the reliability of volunteer contributions. To do that, we compare the work carried out by volunteers and a translation professional in terms of evaluation results and evaluator agreement and examine the control measures used to ensure reliability. Results show similar behaviour regarding general quality assessment but underscore the need for more informative working environments to make evaluation processes more reliable as well as the need for carefully crafted control cases.

1 Introduction

Particularly since trainable neural automatic metrics took centre stage in the WMT metrics shared task in 2022 (Freitag et al., 2022) machine translation (MT) evaluation data sets annotated for quality are becoming essential to develop accurate models. If availing of parallel data with professional references was not difficult enough, we are currently faced with the need to collect data that is

not otherwise produced, in other words, while parallel data could be gathered from previously published translations, a sentence-(or text-)level numeric quality assessment most usually needs to be generated for the specific task of metric training.

This situation poses a particular challenge for low-resource languages which widens the gap between high- and low-resource scenarios. Firstly, because pre-trained models such as COMET (Rei et al., 2020) might not include language-specific data for small languages and therefore quality predictions can be unreliable, and secondly, because collecting relevant annotated data requires a heavy investment. In this context, resorting to opportunistic data collections with crowd volunteers is increasingly tempting.

This new scenario is yet an added reason to increase research efforts on evaluation design. More rigorous considerations of evaluation methodologies and design decisions emerged with claims of human and super-human parity of MT performance (Hassan et al., 2018; Barrault et al., 2019). Researchers claimed that evaluations were not rigorous and pointed out issues such as raters' lack of translation expertise, the quality of reference translations, target language interference in source sentences and non-contextualised evaluations as aspects that skewed results in favour of MT contenders (Läubli et al., 2018; Toral et al., 2018).

Reports of large evaluation initiatives and third-party reviews have shown that little by little evaluation approaches take into account some considerations (Toral, 2020; Popel et al., 2020) and reference campaign such as the annual WMT share task have taken steps to follow best practices for reliable evaluations (Kocmi et al., 2023). Adding to this, research on design-related topics are emerging, such as error methods to opti-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

mise test set configuration to reduce evaluation effort (Saldías Fuentes et al., 2022), classification schemes adapted to identifying critical errors in neural MT (Sudoh et al., 2021), document-level and context-aware agreement and effort (Castilho et al., 2020; Castilho, 2021), detection of post-edited reference translations (Kloudová et al., 2021) and differences between expert and non-expert evaluators (Graham et al., 2013; Freitag et al., 2021). Several crowd evaluation initiatives have also been reported over the years (Bentivogli et al., 2011; Graham et al., 2017) even for low-resource languages (Aranberri et al., 2017; Toral et al., 2017) that cover a number of design decisions. And yet, best practice and efficiency recommendation guidelines are limited and it is not uncommon that evaluation initiatives specially for low-resource scenarios lack the rigour that would benefit the outcomes the most. In this context, the current analysis is only a small step towards studying the characteristics of crowd-based evaluations within minority language communities.

The remaining of this paper is divided as follows Section 2 provides a brief description of the evaluation set-up from where the data set under study originated together with the obtained results; Section 3 outlines the specific details of the evaluation set-up used to obtain a professional evaluation of the said set as well as the qualitative feedback collected on the task and reports a comparison of evaluation results and the agreement analysis between the crowd volunteers and the professional evaluator; Section 4 examines the reliability of the control measures included in the set to identify outlier evaluators; finally, Section 5 draws a number of conclusions from the study.

2 Description of the Original Evaluation Initiative

The data set studied in this work is the product of an evaluation initiative to obtain human assessments of MT for two low-resource languages, namely, Basque and Maltese (Falcão et al., 2024). The authors aimed to collect sentence-level direct assessments to test the potential improvement of the trainable COMET metric with language-specific data. The resulting data set for the Spanish–Basque pair was kindly made available by the researchers for further analysis.¹

¹Access to the data set will be open upon publication of their work.

In this section, we briefly describe the evaluation setup used by the original research (for further details, see Falcão et al. (2024)) and report the overall results for later comparison.

2.1 Evaluation Set-up

Dataset: The evaluation set prepared for the campaign consisted of 400 Spanish source sentences and Basque translations. They were extracted from various existing sets and sources such as FLORES-200², TED2020 (Reimers and Gurevych, 2020), OpenSubtitles (Lison and Tiedemann, 2016), the Elhuyar Corpus³ and the HAC parallel corpus⁴, which cover text from web articles to subtitles and literature. Note that the Spanish source sentences in these sets can include both original and translated text.

Translation sources: The Basque translations paired with the Spanish sentences were obtained from multiple sources. Three MT systems were used to translate the set automatically. Additionally, damaged translations –MT system outputs with an embedded Spanish sequence of words- and reference translations –obtained from the parallel data sets– were also included in the final set as a means to identify unreliable evaluators.

Task: Distributed through the Appraise⁵ platform (Federmann, 2012), the task involved evaluators assessing the translation quality in a continuous scale of 0 to 100. Directed towards a non-specialist participant profile, the description of the task highlighted a series of attributes, including meaning, information, clarity, correctness, grammaticality, and naturalness.⁶ It could be argued

²<https://github.com/facebookresearch/flores/blob/main/flores200/README.md>

³<https://elhuyar.eus/en/services/language-services-and-basque-plan/translation-and-language-resources/corpus>

⁴<https://www.ehu.eus/ehg/hac>

⁵<https://github.com/cfedermann/Appraise/>
<https://github.com/AppraiseDev/Appraise>

⁶The original English text provided in the platform in the relevant languages was as follows: “For each item, you will be shown an original sentence in Spanish and a translation candidate in Basque. You will then be asked to rate the quality of the translation on a scale of 0 to 100, based on how well you believe the translation expresses the full meaning of the original sentence. A rating of 100 means that the candidate is a perfect translation: it expresses the same thing as the original sentence, in a clear and correct manner. A candidate should be rated lower if it contains grammatical or orthographic errors, if it’s missing information, if it sounds unnatural or weird, and so on.” (personal communication, J. Falcão, July 2023)

that the explanation aimed for a general definition of quality rather than a specific aspect. To perform the assessment, evaluators were provided with a source sentence and its corresponding translation. The sentences were provided without context. Participants were free to annotate as many pairs as they wished. No further guidelines were provided as to how to perform the task.

Evaluators: Crowd volunteers were sought by promoting the initiative through university and translator distribution lists, and social media. Therefore, the linguistic profiles of potential evaluators ranges from professional translators to general users with no dedicated training in languages. The evaluators were asked to report their Spanish and Basque language competence to exclude those without an advanced level of both languages. None such cases were reported by the researchers.

2.2 Evaluation Results

A quick analysis of the metadata reveals that 44 crowd volunteers contributed with a total of 1,186 evaluations (an average of 26.95 evaluations per person, with a median of 11). As shown in Table 1, their work is divided as follows: a total of 742 sentence pairs were evaluated,⁷ 389 (%52.42) of which were assessed once, 285 (%38.41) twice and 76 (%10.24) received between three and five annotations. This allowed to collect annotations for about 200 sentences for each MT system (MT1, MT2, MT3), a total of 78 damaged translations, about 25 for each brand of damaged cases (D-MT1, D-MT2, D-MT3), and 53 sentence pairs containing reference translations (Ref).

According to the annotations, MT1 and MT2 score very similarly with results of 77.81 and 78.45 points, respectively (see Table 2).⁸ MT3 lags behind, over 16 points lower. As anticipated, damaged translations score considerably lower, yet following the ranking for the MT systems. Unexpectedly, reference translations score lower than the system averages. As a general trend, the average scores tend to be higher for sentence pairs with a single annotation than for those with multiple annotations. These comparisons should be taken with caution as the sentences annotated for each subgroup are not exactly the same.

⁷Note that this does not cover the over 1,200 pairs in the evaluation set.

⁸Scores for sentence pairs with more than one evaluation were calculated separately; no average was applied.

Trans. source	Sentences	Evaluations			
		1	2	+2	Total
MT1	213	112	78	23	341
MT2	207	97	89	21	342
MT3	191	112	64	15	286
D-MT1	28	14	9	5	48
D-MT2	26	13	10	3	44
D-MT3	24	15	5	4	39
Ref	53	26	22	5	86
Total	742	389	285	76	1186

Table 1: Number of evaluated sentences and collected evaluations, where Trans. source refers to the source from where the translations were obtained, Sentences refers to the number of unique sentence pairs assessed, and evaluations 1, 2 and +2 refer to the number of sentences that obtained the stated number of evaluations.

3 Professional Evaluation

In order to explore the similarity between crowd and professional evaluators and their reliability, the author performed the same evaluation task for the complete set. She is a specialist in translation, native speaker of Basque (accredited C2-level) and Spanish, and has experience in MT evaluation. She will be referred to as the professional evaluator. While the feedback from a single professional is not necessarily indicative of the true annotations the sentence pairs should receive, it can be argued that it provides an educated guess that is consistent across the set to the extent this is possible in human judgement.

3.1 Evaluation Set-up

As in the original evaluation, the evaluator was presented with source and translation pairs to assess in a range of 0-100. The evaluation set consisted of the sentence pairs annotated by the crowd volunteers and 40 additional repeated segments to account for intra-evaluator reliability. The sentence pairs were randomly ordered in a spreadsheet to avoid potential bias and with no access to any additional information (translation source, crowd annotations, etc).

3.2 Qualitative Feedback

Before looking at evaluation results, this section outlines several impressions of the evaluator, noted during the task in an additional column of the spreadsheet, which have been further developed at write-up. While most have already been discussed elsewhere in the literature, this is yet another opportunity to underscore the relevance of evaluation design for reliable and sustainable results.

Trans. source	All evaluations		1 evaluation		+ evaluations		2 evaluations	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD
MT1	77.81	23.73	78.93	24.62	77.26	23.32	77.84	23.16
MT2	78.45	23.32	84.83	21.64	75.92	23.52	74.36	23.73
MT3	61.20	29.17	65.35	27.79	58.53	29.79	57.37	29.36
All MT	73.13	26.45	75.97	26.17	71.72	26.49	70.83	46.57
D-MT1	19.89	21.71	23.78	26.87	18.29	19.44	14.39	16.11
D-MT2	21.59	23.91	22.15	23.87	21.35	24.31	21.20	28.31
D-MT3	12.28	18.30	19.93	24.87	7.50	10.72	6.10	7.00
All D-MT	18.20	21.75	19.88	24.69	16.45	20.12	15.50	21.46
Ref	65.99	28.11	75.08	24.59	62.05	28.82	62.25	27.48

Table 2: Overall evaluation results for each translation source (MT systems, damaged outputs and references) reported as quality mean and standard deviation (SD) broken down per number of evaluations collected for each sentence pair (1, 2 or more and only 2).

The effect of (lack of) context. One of the first topics addressed in translator training is text analysis. Numerous scholars have put forward text analysis frameworks to assist translators in this task. To mention an example, the model for translation-oriented text analysis proposed by Nord (1991) establishes that both extratextual factors (author, sender intention, recipient, medium, place and time of production and reception, motive and function) and intratextual factors (subject matter, hierarchy of content and knowledge presuppositions) should be carefully considered as a first step towards drafting a translation proposal. We see, in fact, that a fully developed translation brief involves information that goes beyond providing the surrounding paragraphs or full text where the translated sentence belongs. And it is only after gathering all those details that a translator can make an informed decision on the adequate register, tone, translation strategies, etc. to be used in their target text. The current evaluation set-up presented sentences in isolation (against the recommendation of the latest WMT campaigns, among others). Reportedly, the result of working without context seems to be that the evaluator favours direct translations, which allow to confirm whether all content and nuances of the source are present in the target language, whereas in a contextual evaluation freer translations that move away from the source to display a more natural use of language and better flow of the text would be accepted and even rewarded. This would be possible because the evaluator would be more informed about the importance of the different contents and formal nuances in the sentences. Conversely, without a clear context, these freer translations can appear

less accurate and may receive a lower score. This behaviour can potentially promote target language words and structures that are more similar to the source language while discouraging the use of expressions that are natural and specific to the target language. Yet another issue brought by the lack of context seems to be that there are cases where it is not easy to judge the correctness of a translation because of ambiguity or incomplete syntactic structure of the source (that is complemented with a previous or following sentence).

The effect of the source. Aggravated in cases where no context is provided and when non-professionals are involved, the source sentence can become somewhat too referential as to what the best translation would be, and might have an impact on scoring, with the evaluator unfairly supporting close wordings (that are grammatical) while undermining more open possibilities that might be more natural and align better with the tone, register and information flow of a text. This can be of particular interest in language pairs for which language contact –and interference into the minority language– is strong and where the vast majority of speakers of the target minority language are also native speakers of the hegemonic language. A (grammar permitting) word for word translation not displaying any target-language specific expressions and structures could be consistently assessed as excellent translations.

The effect of fluency. Accounting for content transfer in translation can be challenging when sentences use complex structures and when subject-knowledge is needed to fully understand the meaning of the source. In these cases, a fluent

translation can be misleading, as extra care is necessary to ensure that all the intended information is present and that no sequences are erroneously interpreted or omitted. This raises the issue of the complexity and thematic typology different evaluator profiles can adequately address.

The precision of the evaluation scale. As a first impression, a 0-100 range seemed very hard to use in the sense that it provided the opportunity to assess quality at a very fine-grained level, while the extent to which mistakes should be penalised seemed greatly subjective. There was a feeling that being consistent with penalisations across the whole evaluation set was hard (see Section 3.3 for agreement results). Admittedly, the evaluator felt more confident with the scale as the number of evaluations performed increased. However, for volunteer work where not a large amount of responses are expected from each individual, a 100-point scale might be a somewhat overwhelming. Note that the professional evaluator wrote a number within an spreadsheet while crowd workers could slide the cursor along a bar, and this might have an impact as well. Additionally, it remained unclear whether the range should be taken as a continuum or a pass/fail threshold should also be considered at 50 points. The annotations collected without the consideration that a score below 50 means that, for example, the translation is unacceptable in a particular situation might differ from those where no such abrupt distinction is made. A similar effect might emerge from scales that use named categories or milestones.

The severity of penalisations. The evaluator reported on the challenge of deciding on a fair penalisation for mistakes. Are 5 points a fair penalisation for an incorrect declension mark? Or should it be 10? 20? 50? Of course, it should depend on the impact it has on the transfer of meaning and on the effect on the form. Even with a context, this is not easy to judge. In fact, impressions noted during the evaluation include a reference to the fact that, not having anchor points to judge the impact of the mistakes and depending on the sentence pair, it would be possible to argue for a score 20 points higher or lower than the one assigned. Moreover, it is not clear whether a penalty should be applied per identified mistake or whether the assigned score should be based on the general impression of the translation quality. To bring a couple of partic-

ularly challenging examples, let us consider the cases where incorrect words or short expressions are encountered that do not align with the overall (good) quality of the rest of the sentence; or cases where a fluent translation that does not convey the same meaning of the original but parts of it are completely correct. Unless the texts are for a specific internal use, sentences with any type of error would most probably be deemed unacceptable. They do not fulfil the intended function of the text. As such, a sentence with a meaning or grammar issue would hardly score above a *pass* threshold in a professional setting. However, when presented with a sentence and a 100-point scale, one might one might penalise a mistake with several points but still assign it a good overall *pass*. Sentences might be evaluated in chunks rather than as a unit.

The effect of the perceived translation competence of evaluators. An idea that emerged during assessment was the extent to which the evaluator's perceived language capacity and translation skills could affect the scoring, in other words, whether evaluators project their own translation competence against the translation provided for assessment and judge according to a self-centred threshold. An evaluator might consider any translation that closely approaches to the quality of what they would produce as a good (or not) translation and score consequently; anything better would be highly valued and weaker sequences of that personal threshold penalised. If this risk exists, it might be pertinent to collect information on the self-perceived translation competence (or actual experience) together with linguistic knowledge.

The post-editing effect. Linked to the issue of error severity is the reported temptation to be more lenient towards important mistakes that are easily fixed and to not assign them a heavy penalisation. The incorrect use of a noun or a preposition that changes the meaning of the whole sentence, for example, but can effortlessly be substituted by an adequate noun or preposition without having to tinker with the rest of the sentence elements can feel less damaging. However, in terms of translation quality, the impact of that incorrect element is crucial. If the aim is to collect quality information, ensuring that evaluators can clearly distinguish between translation adequacy and post-editing effort might be relevant.

The quality of the source segments. During the evaluation task, a considerable amount of source sentences was flagged as including grammar or spelling mistakes. While some did not hinder comprehension, others could obscure the correct interpretation of the intended message. An evaluator will not be able to adequately assess the translation quality of a source sentence they cannot understand. For those sentences that could be (adequately) interpreted despite the mistakes, some translations showed no trace of irregularities and were properly resolved. However, at times, the translations do presents mistakes. The question here is whether we want to penalise a translator’s inability to overcome issues in the source. The presence of problematic source sentences raises the question of the importance of the configuration of the training or evaluation set we aim to gather. It will probably be a good idea to consider the scale of the evaluation (how much data we can collect) and the specific definition of quality we seek and consciously decide whether we want to include not only correct source sentences but also incorrect ones and even variations and levels of *well-writtenness*.

3.3 Evaluation Results

The evaluator assessed a total of 782 segments. The intraclass correlation coefficient (ICC)⁹ calculated with the repeated segments is 0.896 (95% upper bound 0.803 and 95% upper bound 0.9476) which we can interpret as (almost) excellent internal agreement. We visualize these results in an Bland-Altman plot (see Figure 1), where agreement is represented based on the mean difference and by depicting the limits of agreement (Altman and Bland, 1983). If we consider the bias, on average, the second rating of the segments is 2.275 lower, which can be interpreted as a small difference. The data points appear scattered across the graph, indicating the absence of proportional biases or heteroscedasticity. It is also important to note that the great majority of points fall within the limits for the 95% confidence interval, which indicates that the evaluator performs almost equally with the repeated segments. These results can be taken as an indication of consistency in the assessment across the data set. Based on this, we could conclude that the evaluator was able to remain con-

⁹Calculated using a two-way mixed model for absolute agreement for a 95% confidence interval).

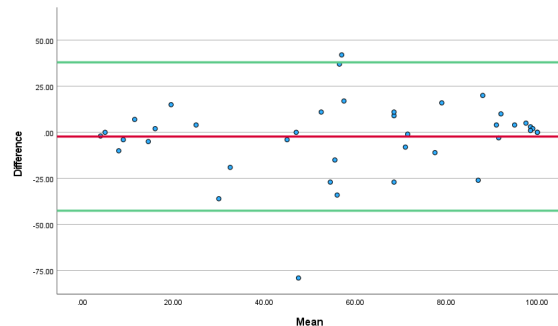


Figure 1: Bland-Altman plot for intra-evaluator agreement for the professional evaluator where the difference between the first and second annotation of the repeated segments is displayed in the Y axis and the average of both annotations is represented in the X axis.

Trans. source	Sent.	Mean	SD	Min.	Max.
MT1	222	78.05	23.07	8	100
MT2	213	79.89	20.94	13	100
MT3	201	58.73	25.26	1	100
All MT	636	72.56	24.94	1	100
D-MT1	32	9.94	8.99	1	30
D-MT2	30	11.07	10.20	1	38
D-MT3	25	8.80	8.71	1	35
All D-MT	87	10.00	9.28	1	38
Ref	59	84.07	18.43	35	100
Total	782	66.47	30.81	1	100

Table 3: Overall evaluation results for each translation source reported as quality mean, standard deviation (SD), minimum score and maximum score according to the annotations of the professional evaluator.

sistent despite the difficulties encountered in the evaluation task and the subjectivity involved in it as described in Section 3.2.

Evaluation results are displayed in Table 3. We can observe that the average score for each of the MT systems is very similar to the scores obtained from the crowd volunteers. These results seem to indicate that overall system quality results would be very similar when evaluated by a translation professional and by (our particular pool of) crowd participants. This is an interesting outcome that might be worth exploring in other evaluation initiatives, as it might be particularly relevant for low-resource scenarios where no funding or professional resources are available for evaluation. It is worth noting that the standard deviations, while large, are slightly smaller than those registered for crowd volunteers and fall within the perceived range of *potential variation* reported by the evaluator (see Section 3.2).

The results for the damaged sentences, however, are up to 10 points lower and all three revolve

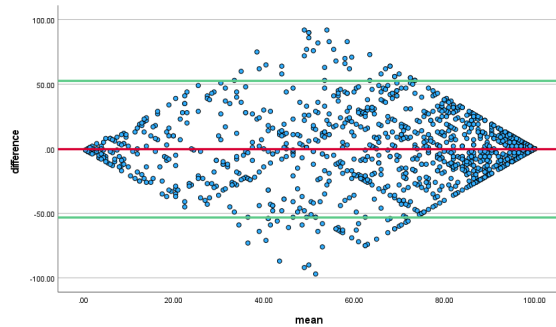


Figure 2: Bland-Altman plot for inter-evaluator agreement between the crowd volunteers and the professional evaluator where the difference between their annotations is displayed in the Y axis and the average of both annotations is represented in the X axis.

around 10 points. Interestingly, even in this case, the ranking for the different damaged sentences follows that of the MT systems that were used to create them. Clearly, the professional translator penalised these cases harsher than volunteers. In contrast, overall, reference translations were assessed almost 20 points higher by the professional evaluator. The differences in these two types of control sentences might indicate that the professional was better equipped to identify the extreme cases and judge them accordingly (see Section 4 for a more thorough analysis of control cases).

In addition to examining quality evaluation results, we also explore the agreement between crowd volunteers and the professional evaluator with respect annotations. Considering all annotations in the data set, the total ICC score is 0.768 (95% lower bound 0.741; 95% upper bound 0.792), which indicates a good agreement (see, for example, Koo and Li (2016) for ICC interpretation). Figure 2 shows a Bland-Altman plot to visualise the overall agreement. The mean difference bias is very close to zero at -0.3125 and we see a random scatter around the mean, mostly within the 95% confidence interval limits. This indicates that the evaluations provided by the two methods observed, that is, a mix of crowd volunteers and a professional evaluator are similar.

If we look more closely, we see that out of the 44 crowd volunteers, when compared with the professional evaluator, three can be assigned an ICC score below 0.5 (poor agreement), 12 an score between 0.5 and 0.75 (moderate agreement), 19 a score between 0.75 and 0.9 (good agreement), and nine a score above 0.9 (excellent agreement).¹⁰ We

¹⁰One evaluator only contributed with one evaluation and was

Agreement level	ICC	95% lower bound	95% upper bound
Poor	-0.044	-0.729	0.381
Moderate	0.700	0.633	0.754
Good and excellent	0.829	0.804	0.851

Table 4: ICC inter-evaluator agreement between the crowd volunteers and the professional evaluator grouped according to the agreement obtained individually.

calculated the ICC scores for the professional evaluator and the groups of crowd volunteers based on the individual level of agreement obtained. The results show that the ICC agreement with the three evaluators with whom the agreement was poor is actually remarkably poor (-0.044), the ICC agreement with those within the moderate range is rather high within that range (0.7) and the ICC agreement with those within the good and excellent range is very good reaching a 0.829 (see Table 4).

If we consider the individual Bland-Altman plots for each subgroup (Figures 3, 4 and 5), we observe that the bias is moving away from zero as evaluators with a lower agreement are represented in the Figures. The scatter seems to widen when comparing the good and excellent group to the moderate group, but it still shows a random pattern. However, the scatter is clearly not random for the evaluators with a poor ICC agreement.

Table 5 shows the evaluation results in terms of translation quality for the crowd volunteers and the professional translator according to the ICC agreement level groups. We can observe that the average quality assigned by agreeing volunteers is similar, whereas the average of the volunteers which agree poorly with the professional differs in over 12 points. However, what is interesting is that across the groups, the stronger disagreements appear for damaged and reference translations, that is, sentences introduced as control elements to identify evaluator reliability. Differences in MT translations remain the most similar and only appear occasionally as agreement levels decrease. We will consider the performance of these control sentences in more detail in Section 4.

4 Discarding Participants

When evaluators are not asked to work on a minimum number of sentences, it becomes highly challenging to identify outliers because it is not possible to perform consistent comparisons. As an approximation to uncover unreliable participants,

therefore not possible to calculate an individual agreement score.

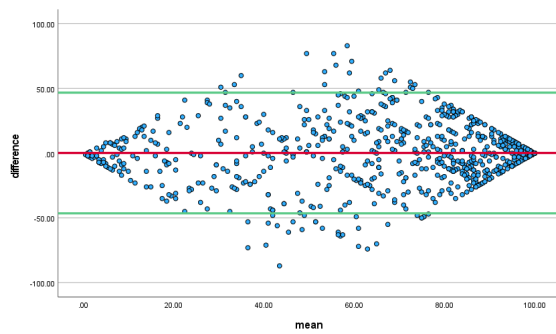


Figure 3: Bland-Altman plot for inter-evaluator agreement between the crowd volunteers and the professional evaluator for which good and excellent ICC agreements were obtained individually, where the difference between their annotations is displayed in the Y axis and the average of both annotations is represented in the X axis.

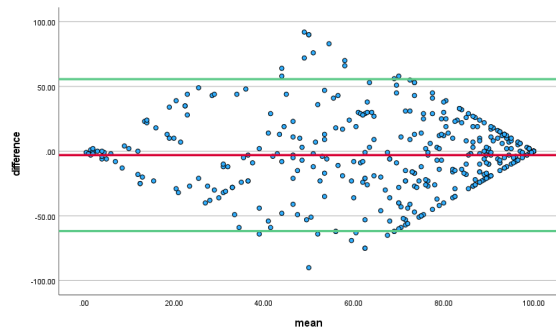


Figure 4: Bland-Altman plot for inter-evaluator agreement between the crowd volunteers and the professional evaluator for which moderate ICC agreements were obtained individually, where the difference between their annotations is displayed in the Y axis and the average of both annotations is represented in the X axis.

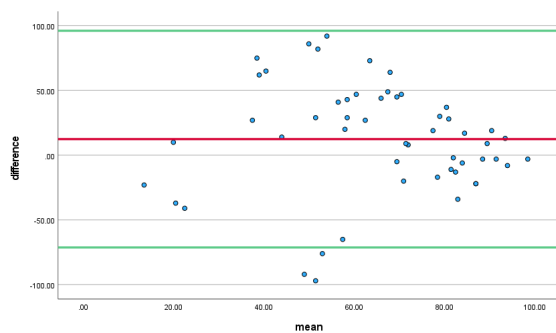


Figure 5: Bland-Altman plot for inter-evaluator agreement between the crowd volunteers and the professional evaluator for which poor ICC agreements were obtained individually, where the difference between their annotations is displayed in the Y axis and the average of both annotations is represented in the X axis.

some works have included in the evaluation set sentence pairs for which the quality is known. It usually involves pairing source segments with artificially damaged translations as examples of bad quality output which should be assessed low, and it can also include source sentences paired with their reference translation as examples of excellent quality pairs. Evaluators who fail to assess them as poor and good cases within a reasonable range are removed from the task and their contributions excluded from the collection. As described in Section 2.2, this is precisely the approach taken by the researchers when running the evaluation task for the data set under study. We looked into these cases to check if these so-called control sentences allow to identify the crowd volunteers which agreed poorly with the professional evaluator.

Out of the 131 damaged sentences included in the data set, 18 were evaluated with a score of 50 or above by 8 crowd evaluators. Out of those 8 evaluators, three have an ICC score ranging between 0.75 and 0.80, that is, a good agreement with the professional evaluator; four have an ICC between 0.61 and 0.72, a moderate agreement; and one of them a low agreement of 0.085, the lowest of all participants. Two out of the three evaluators with a poor ICC agreement with the professional evaluator did not assess any damaged translations. The one who did scored both cases shown with scores above 70 points. This might mean that the inclusion of damaged sentences and its implementation was not particularly accurate in this specific set to filter out deviant evaluators.

To start, not all evaluators assessed damaged sentence pairs. Out of the 44 participants, 36 assessed at least one case. Eight were presented with one damaged sentence pair, 28 were presented with 2 to 11 damaged sentence pairs (Pearson correlation between the total number of evaluations and damaged sentence evaluations is 0.95). Out of them eight failed to pinpoint them as bad quality translations. Only five evaluators with multiple damaged sentence pairs scored them 50 or above more than once. The professional evaluator assigned a low score to all the damaged sentences with scores ranging between 1 and 38.

If we were to discard crowd annotations based on the assessments of damaged sentences, we would be discarding 373 annotations. If we decided to exclude the work of the 8 evaluators who

Trans. source	poor ICC					moderate ICC					good and excellent ICC				
	N	crowd		professional		N	crowd		professional		N	crowd		professional	
		Mean	SD	Mean	SD		Mean	SD	Mean	SD		Mean	SD	Mean	SD
MT1	17	79.47	26.77	71.82	29.64	112	74.62	25.37	78.94	21.18	232	79.07	23.17	77.33	22.86
MT2	15	64.53	31.94	71.93	22.32	108	71.93	25.95	77.93	22.70	228	81.13	22.18	81.62	19.92
MT3	17	66.35	30.91	42.82	24.94	93	62.82	28.57	63.73	25.73	193	59.42	29.42	57.66	24.89
All MT	49	70.35	30.01	61.80	28.99	313	70.74	26.922	74.02	24.02	653	73.98	26.58	73.01	24.66
D-MT1	0	–	–	–	–	21	22.48	30.23	9.33	8.53	31	19.32	18.21	10.55	9.80
D-MT2	3	86.33	11.55	8.67	2.08	10	28.30	22.91	8.00	11.26	37	15.24	18.10	14.68	12.88
D-MT3	0	–	–	–	–	17	19.47	26.25	9.76	10.03	23	8.87	10.60	6.48	5.97
All D-MT	3	86.33	11.55	8.67	2.08	48	22.63	27.10	9.21	9.48	91	15.02	16.885	11.20	10.57
Ref	1	76.00	–	68.00	–	24	58.96	29.84	83.79	16.19	68	71.00	26.98	83.06	19.19
Total	53	71.36	29.17	58.91	30.52	385	63.55	31.28	66.59	31.164	812	67.13	31.676	66.93	30.52

Table 5: Overall evaluation results for each translation source reported as quality mean and standard deviation (SD) for crowd volunteers and the professional evaluator for evaluator groups based on ICC agreements.

were not presented with any damaged translations, we would have to remove another 47 evaluations. This would leave us with a total of 766 evaluations, 64.59% of the total collected. If the damaged translations would have served to accurately identify the outliers (based on the ICC score), we would have only discarded 49 and 47 evaluations, 8.1% of the total evaluations collected.

Let us briefly consider the approach used to create the damaged translations. According to the researchers, these were obtained by translating the source sentences with the three MT systems used to create the remaining data set and by replacing a random sequence of words with a sequence of another source sentence. This approach can result in different types of output, from poor quality target sequences mixed with source language sequences to very good quality target sequences mixed with the source language (note that the quality of the MT systems has been rated within an overall range of 61–78 points). Evaluators trained in translation might be more aware of the importance of the text as a unit and clearly see that such sentences would be unacceptable for the great majority of contexts. Yet this might not be the case for people without translation training. Without a context to consider, it is possible that evaluators do not just penalise the translations for the presence of the source but also feel that they should provide positive points for the sequences with a good quality translation. Depending on the length of the sentences and the proportion of source words, their location within the sentence and the amount of meaning contained in the correct target sequences, it is possible that some evaluators feel that they are being fair by providing a score above 50 to those translations even when they are fully aware of the truncated sequences. Overall, the different results gathered

for damaged sentences in this study might indicate that their current design is probably not the most favourable to serve as reliable control sentences.

Together with damaged translations, reference translations were also included in the data set as a control measure. In total, 86 sentence pairs with references were evaluated. Out of the 44 evaluators, 21 were presented with this type of translations: six assessed one case and the remaining 15 assessed from two to 10 cases. Out of the 86 reference sentences, 26 were evaluated with a score of 50 or below by 13 crowd evaluators. Eight of them, which evaluated two or more of such cases, only assigned this score once, whereas the remaining five assigned a low score in multiple occasions. The professional evaluator assessed 11 of the 86 translations with a score of 50 or lower.

Out of those 13 evaluators, three have an ICC score above 0.9, that is, an excellent agreement with the professional evaluator; 11 have an ICC between 0.76 and 0.87, a good agreement; six have an ICC between 0.61 and 0.74, a moderate agreement, and one of them a poor agreement of 0.085, the lowest of all participants. Once again, two out of the three evaluators with a poor ICC correlation with the professional evaluator did not assess any reference translations. The one who did scored the single case presented with a good score of 76, passing the test. This means that the inclusion of reference sentences was not useful in this specific set to filter out unreliable volunteers, on the contrary, by following this test, we would discard good annotations and keep outlier contributions.

As was the case with damaged translations, not all evaluators assessed reference sentences. Out of the 44 participants, 21 assessed at least one case. Six were presented with one reference sentence pair, 15 were presented with two to ten (Pearson

correlation between the total number of evaluations and damaged sentence evaluations is 0.93). Out of them 13 failed to pinpoint them as good quality translations. Up to five evaluators with multiple reference sentence pairs scored them 50 or below more than once.

If we were to discard crowd annotations based on the assessments of reference sentences, we would be discarding 632 annotations. If we decided to exclude the work of the 23 evaluators who were not presented with any reference translations, we would have to remove another 374 evaluations. This would leave us with a total of 180 evaluations, 15.18% of the total collected. If the reference translations would have served to accurately identify the outliers (based on the ICC score), we would have only discarded 49 and 47 evaluations, 8.1% of the total evaluations collected.

If we take both control measures into account and combine the performance information of the crowd volunteers, we can account for 37 evaluators out of the 44. If we exclude the work carried out by those who failed any of the tests, we would have to remove the contribution of 18 evaluators, that is, a total of 837 evaluations.

Let us briefly consider the case of reference translations. They were extracted from established sets or other bilingual data published as parallel corpora (Falcão et al., 2024). The quality of reference translations included in test sets has often been questioned and so this was investigated further. If we consider the assessment of the professional evaluator, we see that a score of 50 or below was assigned to seven reference translations out of the 59 presented with scores ranging between 37 and 49. The scores assigned to the remaining references varied from 69 to 100. The range is even wider for crowd volunteers, between 51 and 100. This can be a clear indication that the quality of the reference translations was either not always at a professional level or could not be judged as such out of context. Again, this might mean that carefully choosing high quality references and an evaluation set-up that allows to properly assess their quality is important in order to implement an efficient control measure for non-professional initiatives in particular.

5 Final Remarks

In this work we have explored an opportunistic evaluation initiative that targeted a low-resource

language pair (Spanish–Basque) to study the impact of design decisions and the reliability of volunteer participants. A translation professional performed the same evaluation task carried out by volunteers. Next, evaluation results and agreements were compared and the role of control measures that ensure evaluator reliability analysed.

For the analysed set, we can conclude that the overall quality assigned to a MT system might not vary considerably when evaluated by crowd volunteers or a professional evaluator. It remains to be tested if sentence-level accuracy is also as reliable.

In terms of task design, we gathered several issues to consider. Task design is key in providing a working environment that will allow the evaluator to reduce the level of subjectivity and increase consistency. The feedback from a professional evaluator pointed at the benefit of (highly) contextualised sentences, meaningful evaluation categories, manageable complexity and topic specialisation, translation awareness and source sentence quality.

In the same line, identifying outlier contributions seems key to guaranteeing a reliable annotated data set. This being the case, our analysis has demonstrated that while damaged and reference sentences might serve as measures to identify unreliable participants, attention must be paid to creating them. The resulting translations must be unquestionably poor/good so that alternative interpretations are ruled out. Then again, it remains to be studied whether participants who properly assess control sentences that are too easily identifiable as poor/good translations will be able to accurately assess regular MT output quality.

All in all, we must not forget that these remarks emerge from the analysis of a single data set and a particular crowd volunteer group. In fact, it would be interesting to study if there are commonalities among the characteristics of the crowd volunteer communities of minoritised languages (participant profiles, level of commitment, level of agreement with professional assessment, for example) and whether these are similar to the crowd participants of hegemonic languages.

Also, this work has explored design issues that are relevant in terms of translation assessment and reliability. However, further research into the real impact of more accurate and cleaner annotations on model training would also be beneficial to determine how rigid (or flexible) an evaluation set-up must be in order to yield useful annotations.

Acknowledgements

This work is supported by the TRAIN (PID2021-123988OB-C31) project funded by MCIN/AEI/ 10.13039/501100011033 and by “ERDF A way of making Europe”, DeepR3 (TED2021-130295B-C31) founded by MCIN/AEI/10.13039/501100011033 and European Union NextGeneration EU/PRTR, and the Basque Government (IXA excellence research group IT1570-22).

References

- Altman, Douglas G and J Martin Bland. 1983. Measurement in medicine: the analysis of method comparison studies. *Journal of the Royal Statistical Society Series D: The Statistician*, 32(3):307–317.
- Aranberri, Nora, Gorka Labaka, Arantza Díaz de Ilarraza, and Kepa Sarasola. 2017. Ebalua-toia: crowd evaluation for english–basque machine translation. *Language Resources and Evaluation*, 51:1053–1084.
- Barrault, Loïc, Ondřej Bojar, Marta R Costa-Jussa, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, et al. 2019. Findings of the 2019 conference on machine translation (wmt19). ACL.
- Bentivogli, Luisa, Marcello Federico, Giovanni Moretti, and Michael Paul. 2011. Getting expert quality from the crowd for machine translation evaluation. In *Proceedings of Machine Translation Summit XIII: Papers*, Xiamen, China, September 19-23.
- Castilho, Sheila, Maja Popović, and Andy Way. 2020. On context span needed for machine translation evaluation. In Calzolari, Nicoletta, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asuncion Moreno, Jan Odiijk, and Stelios Piperidis, editors, *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3735–3742, Marseille, France, May. European Language Resources Association.
- Castilho, Sheila. 2021. Towards document-level human MT evaluation: On the issues of annotator agreement, effort and misevaluation. In Belz, Anya, Shubham Agarwal, Yvette Graham, Ehud Reiter, and Anastasia Shimorina, editors, *Proceedings of the Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 34–45, Online, April. Association for Computational Linguistics.
- Falcão, Júlia, Kurt Arbela, Nora Aranberri, and Claudia Borg. 2024. Comet for low-resource machine translation evaluation: A case study of english-maltese and english-basque. In *LREC-COLING 2024 - The 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation*, Turin, Italy.
- Federmann, Christian. 2012. Appraise: an open-source toolkit for manual evaluation of mt output. *Prague Bull. Math. Linguistics*, 98:25–36.
- Freitag, Markus, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.
- Freitag, Markus, Ricardo Rei, Nitika Mathur, Chikui Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Graham, Yvette, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. Continuous measurement scales in human evaluation of machine translation. In Pareja-Lora, Antonio, Maria Liakata, and Stefanie Dipper, editors, *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 33–41, Sofia, Bulgaria, August. Association for Computational Linguistics.
- Graham, Yvette, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2017. Can machine translation systems be evaluated by the crowd alone. *Natural Language Engineering*, 23(1):3–30.
- Hassan, Hany, Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, et al. 2018. Achieving human parity on automatic chinese to english news translation. *arXiv preprint arXiv:1803.05567*.
- Kludová, Věra, Ondřej Bojar, and Martin Popel. 2021. Detecting post-edited references and their effect on human evaluation. In Belz, Anya, Shubham Agarwal, Yvette Graham, Ehud Reiter, and Anastasia Shimorina, editors, *Proceedings of the Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 114–119, Online, April. Association for Computational Linguistics.

- Kocmi, Tom, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. 2023. Findings of the 2023 conference on machine translation (wmt23): Lms are here but not quite there yet. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42.
- Koo, Terry K and Mae Y Li. 2016. A guideline of selecting and reporting intraclass correlation coefficients for reliability research. *Journal of chiropractic medicine*, 15(2):155–163.
- Läubli, Samuel, Rico Sennrich, and Martin Volk. 2018. Has machine translation achieved human parity? a case for document-level evaluation. In Riloff, Ellen, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4791–4796, Brussels, Belgium, October–November. Association for Computational Linguistics.
- Lison, Pierre and Jörg Tiedemann. 2016. Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles.
- Nord, Christiane. 1991. *Text Analysis in Translation*. Rodopi, Amsterdam, Atlanta.
- Popel, Martin, Marketa Tomkova, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar, and Zdeněk Žabokrtský. 2020. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. *Nature communications*, 11(1):1–15.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In Webber, Bonnie, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.
- Reimers, Nils and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation. In Webber, Bonnie, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4512–4525, Online, November. Association for Computational Linguistics.
- Saldías Fuentes, Belén, George Foster, Markus Freitag, and Qijun Tan. 2022. Toward more effective human evaluation for machine translation. In Belz, Anya, Maja Popović, Ehud Reiter, and Anastasia Shimorina, editors, *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 76–89, Dublin, Ireland, May. Association for Computational Linguistics.
- Sudoh, Katsuhito, Kosuke Takahashi, and Satoshi Nakamura. 2021. Is this translation error critical?: Classification-based human and automatic machine translation evaluation focusing on critical errors. In Belz, Anya, Shubham Agarwal, Yvette Graham, Ehud Reiter, and Anastasia Shimorina, editors, *Proceedings of the Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 46–55, Online, April. Association for Computational Linguistics.
- Toral, Antonio, Miquel Esplá-Gomis, Filip Klubička, Nikola Ljubešić, Vassilis Papavassiliou, Prokopis Prokopidis, Raphael Rubino, and Andy Way. 2017. Crawl and crowd to bring machine translation to under-resourced languages. *Language resources and evaluation*, 51:1019–1051.
- Toral, Antonio, Sheila Castilho, Ke Hu, and Andy Way. 2018. Attaining the unattainable? reassessing claims of human parity in neural machine translation. In Bojar, Ondřej, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 113–123, Brussels, Belgium, October. Association for Computational Linguistics.
- Toral, Antonio. 2020. Reassessing claims of human parity and super-human performance in machine translation at WMT 2019. In Martins, André, Helena Moniz, Sara Fumega, Bruno Martins, Fernando Batista, Luisa Coheur, Carla Parra, Isabel Trancoso, Marco Turchi, Arianna Bisazza, Joss Moorkens, Ana Guerberof, Mary Nurminen, Lena Marg, and Mikel L. Forcada, editors, *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 185–194, Lisboa, Portugal, November. European Association for Machine Translation.

Implementations & Case Studies

A Case Study on Contextual Machine Translation in a Professional Scenario of Subtitling

♣♦ Sebastian Vincent, ♦ Charlotte Prescott ♦ Chris Bayliss,
♦ Chris Oakley, ♣ Carolina Scarton

♣ Department of Computer Science, University of Sheffield, UK
♦ ZOO Digital Group PLC, UK

Abstract

Incorporating extra-textual context such as film metadata into the machine translation (MT) pipeline can enhance translation quality, as indicated by automatic evaluation in recent work. However, the positive impact of such systems in industry remains unproven. We report on an industrial case study carried out to investigate the benefit of MT in a professional scenario of translating TV subtitles with a focus on how leveraging extra-textual context impacts post-editing. We found that post-editors marked significantly fewer context-related errors when correcting the outputs of MTCUE, the context-aware model, as opposed to non-contextual models. We also present the results of a survey of the employed post-editors, which highlights contextual inadequacy as a significant gap consistently observed in MT. Our findings strengthen the motivation for further work within fully contextual MT.

1 Introduction

As an innovation-driven company offering dubbing and subtitling services, ZOO Digital is dedicated to exploring assistive technologies to streamline our workflows. Machine translation in particular is a promising tool for improving the efficiency of the (currently fully manual) translation of the transcribed video content during interlingual subtitling. Our domain is characterised by specific challenges, both linguistic (preservation of

style and function in dialogue) and practical (keeping within subtitle constraints, such as visual properties and considerations for the viewers' reading speed). We report on a case study where translation from scratch was replaced with post-editing machine translations of the source text. While such a formulation is far from new – MT has been consistently demonstrated to help reduce effort in the subtitling domain (C. M. de Sousa et al., 2011; Huang and Wang, 2023) – previous studies have relied on off-the-shelf general-purpose neural machine translation (NMT) engines like Google Translate¹. Our work investigates two additional systems: BASE-NMT, a specialised engine trained on our data, as well its contextual version based on the MTCUE architecture (Vincent et al., 2023), whose training involves observing a vast range of metadata and document-level information.

The study was carried out with the assistance of translation and post-editing professionals. Hereinafter we refer as *post-editors (PEs)* to those who were tasked with post-editing work, and as *translators (HTs)* to those who were tasked with translation from scratch (FST). The campaign took place in a full-context multi-modal environment where the professionals had access to the video material and were able to directly jump to the segment corresponding to the utterance they were reviewing, as well as see the preceding and succeeding segments. A total of eight PEs were employed, four for English-to-German (EN-DE) and four for English-to-French (EN-FR) translation, and four HTs, two per language pair. We measured the effort it took to post-edit or translate the TV series content and the number of specific translation errors observed by the PEs. Our findings highlight

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹<https://translate.google.com/>

the necessity of tailoring MT engines to the target domain and motivate further work within leveraging contextual systems in dialogue translation.

2 Related Work

Over the last few years, subtitle translation has been given a volume of attention: C. M. de Sousa et al. (2011), Koponen et al. (2020) and Huang and Wang (2023) observe that post-editing the outputs of an NMT system is a promising alternative to translation *ex novo*, reducing the temporal, technical and cognitive effort of both novice and professional translators and subtitlers. A survey among professional subtitlers detailed by Karakanta et al. (2022), finds that professionals have a positive outlook on incorporating automatic components (such as MT) into their workflow, as they offer starting templates, reduce effort and can provide useful suggestions. However, some challenges in the automatic translation of subtitles remain unsolved (Gupta et al., 2019; Karakanta et al., 2022), including the adherence to subtitle block limitations, which often necessitates shorter and paraphrased translations; lexical consistency, which involves translating the same terms across the text, as well as using vocabulary that maintains the cohesion and coherence of the text, aligns with the surrounding video or textual content, and conforms to standard language or industry conventions; lexical errors such as the translation of idioms and figurative language, and context-related inconsistencies. Context-related errors in particular have been pointed out as the culprit in many works in MT that leveraged the OpenSubtitles corpus (Lison et al., 2018), a dataset of user-submitted subtitles and their translations. Leveraging document-level information (Tiedemann and Scherrer, 2017; Bawden et al., 2018), speaker’s and interlocutor’s gender identity (Vincent et al., 2022) and explicit extra-textual information (Vincent et al., 2023) has been found particularly useful in addressing this challenge. Context is also useful during the manual post-editing procedure: Huang and Wang (2023) show that such a setup decreases the cognitive load of student translators compared to a text-only scenario, suggesting as an explanation the dual coding theory, according to which the interactions between the verbal and non-verbal information enhances the translators’ understanding of the material.

This work employs MTCUE (Vincent et al.,

2023), a multi-encoder Transformer designed for contextual NMT capable of leveraging contextual signals such as film metadata and document-level information to improve translation quality, as well as enabling better control of phenomena such as speaker’s gender and formality register. The mechanism for delivering context in the model involves converting the context fields into equal-sized vectors via sentence embedding. The resulting vector sequence is inputted into a distinct Transformer encoder. Additionally, we employ the context specificity evaluation method outlined in Vincent et al. (2024), which relies on the pointwise mutual information (PMI). In this method, PMI quantifies the degree of co-occurrence between tokens in a translation hypothesis and the respective context.

3 Experimental Setup

The primary objective of our case study was to investigate whether post-editing MT is a cost-effective alternative to FST in our workflow, and to what extent domain-adapted training data and the utilisation of context have an impact in this area. Guided by the availability of resources, we operated in two language pairs: EN-DE and EN-FR and considered four versions of the text in each, including MT outputs from three systems:

1. GOOGLE², a general-purpose NMT engine used in previous work.
2. BASE-NMT, a non-contextual Transformer-based translation model parameter-matched to MTCUE and trained on the same data (except context).
3. MTCUE system (Vincent et al., 2023), a multi-encoder Transformer.

We also operated on the human translations of the test set (REF) approved for production.³ For both MTCUE and BASE-NMT, we trained the models after Vincent et al. (2024), §4.1, in the OVERLAP setting which mimics a scenario with access to prior episodes of a tested series for training (a sample is presented in Appendix F of that work). We operated on sentence-level translations, with MTCUE using the context for each sentence in its dedicated space.

²<https://translate.google.com/>

³This baseline is omitted during automatic evaluation (in fact, it is used as the reference text to calculate the automatic metrics), but is used as a baseline in the human evaluation, where the professionals are asked to post-edit this already sufficiently good text.

3.1 Automatic evaluation

We conducted a pre-emptive automatic evaluation to confirm the feasibility of the human evaluation study. We used BLEU (Papineni et al., 2002) and COMET (Rei et al., 2020) as translation quality metrics. Additionally, to measure context specificity, we measured the PMI between contextual and non-contextual translations (Vincent et al., 2024). We compared the outputs of the machine translation systems (BASE-NMT, GOOGLE, MTCUE) against the reference (REF).

3.2 Post-Editing Setup and Metrics

The human evaluation aspect of the study is interpreted as the effort required to post-edit the translations to a production standard, and captured in the **number of errors, keystrokes and total edit time**. The task was performed by professional HTs and PEs using ZOOSUBS, an in-house software application belonging to ZOO Digital, built to facilitate manual translation of video material (Figure 1). The software’s interface displays the video material along with timed subtitles in the original language. The *target stream*, i.e. the set of text boxes provided to the right of the source stream, is where the HTs input their translations to the desired language. It can optionally be pre-populated with “draft” translations – a setting we opted for in this study – allowing post-editors to edit, divide or combine the segments as they see fit.

To make amendments to a segment, the PE needs to click on its box. From that point, the system tracks the time spent editing the box and the number of keystrokes made. These metrics are recorded for each box separately and taken into account only if the post-edited text differs from the original. After applying modifications, an **Issues for event** window appears for the user to specify the purpose of the changes by selecting errors from a predefined list, optionally providing text commentary. We leveraged this functionality of ZOOSUBS to measure the total and average time and number of keystrokes made by HTs and PEs given some pre-existing translations. We also measured the number of selected errors. For this project, we created a bespoke taxonomy of errors (Table 1) based on translation errors reported in previous work (Freitag et al., 2021; Sharou and Specia, 2022), the original list of issues already present in the ZOOSUBS system and relevant errors from previous work (§2). Error categories

from the aforementioned sources were compiled together and curated to fit the study requirements⁴

Worker setup The PEs operated on seven episodes from three TV series of varying genres: a fictional series about space exploration, a documentary exploring aspects of everyday life, and a family cooking competition show. They were unaware that some of the text they worked with was machine translated, but were told that it was for a research project and asked to relax some constraints such as adhering to the reading speed limits. In addition, we asked four HTs (two to German, two to French) to translate one episode of the cooking show from scratch in ZOOSUBS so we could compare their effort to that of post-editors. For each of the seven episodes, the PEs were asked to post-edit one out of four versions of the text, corresponding to the list outlined in §3. We included the human references (REF) to account for the fact that PEs can sometimes post-edit a translation even when the original one is valid. Our setup ensured that the same PE evaluated the output for each episode exactly once (i.e. does not see two different versions of the same text) (Table 2). When referring to individual PEs, we use the notation $PE.[L][i]$, where $L \in \{G \text{ (German), } F \text{ (French)}\}$, and i denotes the PE ID $\in [1, 4]$.

Details regarding the PEs The recruited PEs and HTs were professionals within the subtitle domain and freelance employees of ZOO DIGITAL. They were informed that the undertaken work was carried out for a research project, but nevertheless, they were paid for their effort at competitive PE and HT rates, standard within the company for this type of work. Information about the PEs’ and HTs’ years of experience (YOE) was collected to shed more light on the findings (Table 3). They also answered a short survey about their views regarding machine translation, discussed in detail in §5.3:

1. Which one would you prefer: translating a stream from scratch or completing a quality check on (post-editing) a stream? Why?
2. What are your views on the use of machine translation in the industry?
3. In your view, are there benefits to post-editing translations over translating from scratch?

⁴We uploaded a draft taxonomy to ZOOSUBS, and the first author performed a test evaluation against a stream with 446 segments to validate the list. As a result, some errors were split into more granular categories, some were renamed and some generalised.

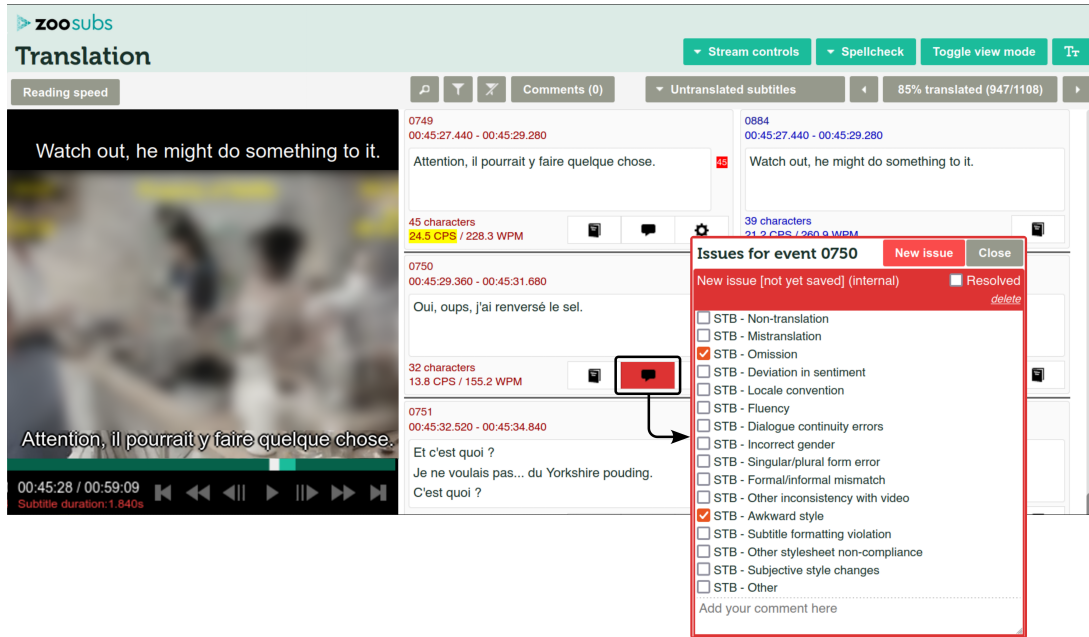


Figure 1: A compressed snapshot of ZOOSUBS.

Type	Description
Translation quality	
<i>Catastrophic translation</i>	Impossible to post-edit, must be translated from scratch.
<i>Mistranslation</i>	Incorrect. Does not preserve the meaning or function of the source.
<i>Omission</i>	Part of the source text was left untranslated.
<i>Deviation in sentiment</i>	Does not preserve the sentiment of the source (e.g. does not match the expressed excitement), or negates the sentiment (e.g. from positive to negative).
<i>Locale convention</i>	Violates locale convention, e.g. currency and date format.
<i>Fluency</i>	Contains punctuation, spelling and grammar errors.
Context	
<i>Incorrect gender</i>	Misgenders the speaker or the addressed person(s).
<i>Incorrect plurality</i>	Incorrectly refers to a single person when a group is addressed, or vice versa.
<i>Wrong formality</i>	Expressed in informal style or uses informal addressing when should use formal, or vice versa.
<i>Other inconsistency with video</i>	Contains inconsistencies with the video material not falling within any of the above.
Style	
<i>Subtitle formatting violation</i>	Violation of the subtitle blocking guidelines.
<i>Other style sheet non-compliance</i>	Does not conform to the provided style sheet.
<i>Awkward style</i>	The style of the translation does not reflect the style of the source sentence and/or the context.
<i>Subjective style changes</i>	The translation is acceptable but the editor suggests improvements in style.
Other	Error of type not found above (use text box provided).

Table 1: List of errors provided to the human evaluators during the campaign.

Series	A		B		C		
	A1	A2	B1	B2	C1	C2	C3
PE.1	REF	MTCUE	GOOGLE	BASE-NMT	REF	MTCUE	GOOGLE
PE.2	BASE-NMT	REF	MTCUE	GOOGLE	BASE-NMT	REF	MTCUE
PE.3	GOOGLE	BASE-NMT	REF	MTCUE	GOOGLE	BASE-NMT	REF
PE.4	MTCUE	GOOGLE	BASE-NMT	REF	MTCUE	GOOGLE	BASE-NMT
HT.1	From Scratch						
HT.2	From Scratch						

Table 2: Work assignment to PEs and HTs in the human evaluation campaign used for both language pairs.

	English-to-French				English-to-German			
	PE.F1	PE.F2	PE.F3	PE.F4	PE.G1	PE.G2	PE.G3	PE.G4
Translation YOE	15	8	3	20	7	18	8	17
YOE in subtitles	8	6	1.5	20	7	5	8	7
YOE in post-editing	8	6	3	10	5	5	1	3
Post-editing training?	✓	✓	✓	✓	✗	✗	✗	✗
Prefer post-editing?	✓	✓	✗	✓	✓/✗	✗	✗	✗

Table 3: Details regarding employed PEs.

All French HTs had training in post-editing, and three out of four preferred it to translating from scratch, while no German HTs had received such

training in the past, and all but one strictly preferred FST. All PEs had at least one YOE in post-editing and one and a half in the subtitle domain. Although the HTs within both pairs had a similar

amount of experience in translation in general and in the subtitle domain (11.5 ± 6.5 for French vs 12.5 ± 5.0 for German), the French HTs had the advantage in terms of YOE in both subtitling (a mean difference of 2.1 YOE) and post-editing (a mean difference of 3.3 YOE).

4 Results of Automatic Evaluation

The automatic evaluation results (Figure 2) suggest that MTCUE was the best-performing system and GOOGLE the worst-performing for both language pairs. Interestingly, for EN-DE, the BLEU and COMET score differences varied in magnitude, to the point of COMET judging all three systems as on par. A possible cause was the discrepancy in hypothesis length (the reference text uses 7.04 words per segment, BASE-NMT: 7.06, MTCUE: 7.06, GOOGLE: 8.29). Since COMET’s calculation involves comparing sentence embeddings of the hypothesis and the reference, including more words or phrases in the hypothesis may lead to a closer similarity match, inflating the score even if the additional tokens are redundant or even harmful to quality. BLEU does not have this problem as it is based on string matching (Papineni et al., 2002). As per the PMI scores, the professional translations (REF) consistently exhibited the highest context specificity. However, MTCUE was on par with this reference score in both cases and was consistently better than the other two systems. MTCUE therefore shows promise at addressing the context-related issues in subtitle translation.

5 Results of the Post-Editing Study

This section analyses the results of the post-editing study: the translation errors (§5.1), the post-editing effort (§5.2), and finally, the post-campaign survey responses (§5.3).

Due to the unprecedented nature of this work in the company, the professionals’ contract allowed them to withdraw if they found the compensation insufficient for the requested work. At the midpoint of the campaign, two PEs (**PE.G1** and **PE.G3**) contacted the project manager to express concerns regarding the quality of the MT outputs, asserting that the task potentially required more effort than FST. To compromise, they proposed narrowing the scope of the remaining work to error identification and marking, without making the necessary corrections. This meant we would not obtain the effort metrics for the two PEs. Conse-

quently, while the error analysis in §5.1 includes both language pairs, the effort analysis in §5.2 does not include results from **PE.G1** or **PE.G3**.

5.1 Error Analysis

An initial inspection of the results indicated that each PE marked a significantly different total number of errors (e.g. **PE.F1** marked 232 errors total while **PE.F4** marked 878). This made direct comparison of the error counts across systems unreliable as each PE also post-edited a different number of segments for each system (cf. Table 2). With seven episodes and four different versions of the text, for each PE there is a version of text they would only have seen one episode from. For example, in Table 2, **PE.1** is assigned two episodes for REF, MTCUE and GOOGLE, but only one for BASE-NMT. In this example, if **PE.1** generally marked fewer errors than others, BASE-NMT would be disproportionately rewarded.

To make the measurements comparable, we normalised them by computing a *normalisation coefficient* h for each PE and then multiplying their error counts for each category by their h . Let $ERR_{PE_i,c}$ denote the number of errors within the category c for the i -th PE. We compute the normalised count $\widehat{ERR}_{PE_i,c}$ as described by Equation 1.

$$\widehat{ERR}_{PE_i,c} = ERR_{PE_i,c} \times h_i$$

$$\text{where } h_i = \frac{\max(ERR_{PE_j,total}; j \in \{1, 4\})}{ERR_{PE_i,total}} \quad (1)$$

We report the total error counts as well as the normalisation multipliers in Table 4.

English-to-German			English-to-French		
PE ID	Error count	h	PE ID	Error count	h
PE.G1	1526	1.76	PE.F1	232	14.68
PE.G2	2452	1.10	PE.F2	182	18.71
PE.G3	2690	1.0	PE.F3	3406	1.0
PE.G4	1832	1.47	PE.F4	878	3.88

Table 4: Error counts and values of h for each PE.

Error post-processing To facilitate post-editing in ZOOSUBS, MT outputs had to be adapted to match the subtitle format. Quality checks of translations conducted in ZOOSUBS normally require the users not just to ensure the correctness of translations but also that the subtitles comply with strict guidelines⁵. Typical MT systems, like the ones

⁵This includes adhering to reading speed and length limits, balancing the length of the top and bottom subtitle, disam-

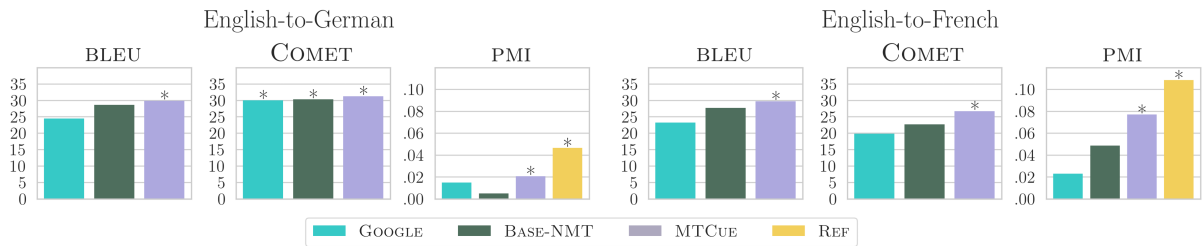


Figure 2: BLEU, COMET and PMI scores obtained by the evaluated models. Asterisks (*) over bars indicate the best result along with all statistically indistinguishable results computed either via bootstrap resampling (or t-test for PMI), $p = 0.05$.

used in this project, are not designed to create translations conforming to these stringent guidelines, and the primary goal of this study was to identify the impact of the translation errors alone. To faithfully replicate the normal work environment of the PEs, we applied a greedy reformatting tool (built into ZOOSUBS) to reformat our translations as subtitles. We made it clear that the project is centred on the correctness of translations, not the subtitle formatting. Still, to ensure that the translation and non-translation errors are kept separate, we included two environment-specific errors for the workers to select from: *Subtitle formatting violation* covering cases where the subtitle is not split to optimally adhere to segmentation guidelines; and *Other style sheet non-compliance* where a rule outlined in the style sheet from the client company was not followed, such as custom punctuation conventions.

Example 1		Target: German
Source	Can I take a look at what you're doing by any chance?	
BASE-NMT (X)	Kann ich mir zufällig ansehen, was du [BR] machst?	
Post-ed.	Kann ich mir vielleicht ansehen, [BR] was Sie da machen?	
Errors	<i>Mistranslation</i> <i>Subtitle formatting violation</i> <i>Formal/informal mismatch</i>	

In some instances, a PE would encounter both translation and non-translation errors within the same segment, as exemplified in **Example 1**, where both translation errors (*Mistranslation of by any chance* and *Formal/informal mismatch of you're doing*) and non-translation errors (*Subtitle formatting violation* of the position of the subtitle break) are present. In such cases, we (i) disregard the non-translation error counts, and (ii) correct

biguation of speaker turns with colours or dashes, and applying appropriate formatting, as specified by a style sheet.

the effort rates (editing time and keystrokes) to account solely for translation-related errors. To precisely gauge the latter, we employed a correction method: let $ERR_{non-translation}$ and $ERR_{translation}$ be the total effort expended by a PE on a segment that had only non-translation and only translation errors marked, respectively. We calculated translation share (TS) as follows:

$$TS = \frac{ERR_{translation}}{ERR_{translation} + ERR_{non-translation}}$$

We then used it to calculate the estimated share of the effort spent on translation in segments that had both errors marked by multiplying TS by the total effort spent on a segment with both error types.⁶

Finally, since the **Other** category was used substantially, we parsed the contents of the optional description text box. The most commonly reported **Other** errors were “Grammar”, “Punctuation”, “Timing”, “SGP” (spelling, grammar, punctuation) and “Literal translation”. Such errors (69.3%) were removed from the **Other** category and pigeonholed as appropriate (e.g. “Grammar” as *Fluency*). More complex comments such as “wissen Sie should not be in the translation” were left categorised as *Other* (30.7%).

Results The calculated normalised counts of errors within each category (Table 5) suggest that MTCUE performs no worse than both non-contextual MT systems overall (row **Total**), while performing significantly better in the **Context** and **Style** categories in EN-FR, pointing to gains related to the use of context information.

The most frequently flagged errors in both language pairs were consistently *Mistranslation* and *Fluency*. *Mistranslation* was reported a similar number of times for all three machine translation

⁶For example, if a PE took three seconds for translation errors and two seconds for non-translation errors on average, where they marked both types we multiplied their total effort for that segment by $\frac{3}{3+2}$.

Error type	Normalised count			
	GOOGLE	BASE-NMT	MTCUE	REF
Translation quality	13.12 ± 14.46	<u>8.70 ± 11.67</u>	8.49 ± 10.90	4.56 ± 5.14
<i>Catastrophic translation</i>	<u>0.50 ± 0.27</u>	0.46 ± 0.18	<u>0.88 ± 0.95</u>	0.72 ± 0.68
<i>Mistranslation</i>	<u>26.99 ± 8.58</u>	25.69 ± 7.67	<u>26.74 ± 6.15</u>	8.76 ± 5.51
<i>Omission</i>	0.26 ± 0.15	2.32 ± 2.20	3.54 ± 2.79	5.38 ± 6.75
<i>Deviation in sentiment</i>	<u>1.11 ± 0.66</u>	0.83 ± 0.30	<u>1.25 ± 0.88</u>	5.23 ± 4.40
<i>Locale convention</i>	2.04 ± 0.00	<u>0.94 ± 0.46</u>	0.61 ± 0.30	0.91 ± 1.03
<i>Fluency</i>	16.88 ± 15.22	<u>9.54 ± 11.17</u>	7.10 ± 6.52	4.18 ± 3.65
Context	5.34 ± 5.68	<u>2.64 ± 3.45</u>	2.21 ± 2.55	1.18 ± 1.13
<i>Incorrect gender</i>	<u>2.20 ± 1.58</u>	<u>1.69 ± 1.90</u>	1.43 ± 1.17	1.60 ± 1.19
<i>Plural/singular form error</i>	<u>0.99 ± 0.81</u>	0.80 ± 0.63	<u>1.19 ± 1.24</u>	0.33 ± 0.00
<i>Formal/informal mismatch</i>	11.31 ± 4.55	<u>5.29 ± 4.60</u>	3.86 ± 3.60	1.19 ± 1.31
Style	<u>12.19 ± 9.79</u>	8.12 ± 6.59	<u>9.88 ± 7.83</u>	3.77 ± 3.86
<i>Awkward style</i>	17.70 ± 7.76	11.82 ± 5.21	<u>13.11 ± 7.04</u>	4.70 ± 4.34
<i>Subjective style changes</i>	<u>2.55 ± 2.09</u>	1.65 ± 1.59	<u>2.33 ± 2.28</u>	2.13 ± 2.52
Other	<u>2.12 ± 3.43</u>	<u>3.26 ± 4.48</u>	2.10 ± 2.46	3.39 ± 5.88
Total	9.58 ± 11.35	<u>6.44 ± 9.05</u>	6.41 ± 8.82	3.86 ± 4.70
Translation quality	20.01 ± 23.05	9.27 ± 9.52	<u>10.21 ± 8.88</u>	6.60 ± 5.08
<i>Catastrophic translation</i>	<u>3.41 ± 1.38</u>	2.25 ± 2.39	<u>2.86 ± 3.03</u>	2.51 ± 3.26
<i>Mistranslation</i>	38.80 ± 14.35	<u>22.73 ± 8.49</u>	20.10 ± 7.34	7.24 ± 3.61
<i>Omission</i>	2.40 ± 2.40	<u>3.91 ± 1.49</u>	5.56 ± 4.09	7.48 ± 5.13
<i>Deviation in sentiment</i>	5.93 ± 5.90	<u>7.82 ± 6.09</u>	11.59 ± 0.00	6.74 ± 3.03
<i>Locale convention</i>	4.29 ± 2.49	0.73 ± 0.51	0.21 ± 0.00	0.63 ± 0.00
<i>Fluency</i>	30.83 ± 31.77	<u>7.28 ± 3.75</u>	5.92 ± 4.18	7.82 ± 7.35
Context	<u>5.41 ± 3.64</u>	6.09 ± 4.26	3.86 ± 3.11	1.29 ± 1.07
<i>Incorrect gender</i>	3.49 ± 2.59	6.96 ± 5.57	<u>4.77 ± 3.98</u>	0.49 ± 0.44
<i>Plural/singular form error</i>	4.50 ± 1.92	5.84 ± 4.60	1.97 ± 0.62	0.00 ± 0.00
<i>Formal/informal mismatch</i>	<u>7.44 ± 4.63</u>	<u>5.58 ± 3.76</u>	4.23 ± 2.93	1.69 ± 1.10
Style	11.05 ± 7.07	10.35 ± 3.69	3.41 ± 2.53	5.55 ± 3.41
<i>Awkward style</i>	11.13 ± 7.46	9.55 ± 1.27	2.89 ± 2.76	4.10 ± 1.28
<i>Subjective style changes</i>	<u>10.94 ± 8.16</u>	11.15 ± 5.52	4.18 ± 2.87	6.28 ± 4.09
Other	37.20 ± 52.68	11.19 ± 16.44	<u>23.67 ± 29.23</u>	27.05 ± 24.68
Total	17.02 ± 25.78	8.84 ± 9.20	<u>9.63 ± 13.85</u>	8.83 ± 12.84

Table 5: Counts of errors flagged by the PEs for each system. Excluding REF, the best result in each row is highlighted and all statistically indistinguishable results are underlined (one-tailed t-test, confidence interval of 80%, $p = 0.2$). Error rates for categories in bold (e.g. **Style**) are calculated based on all errors within the category.

systems in EN-DE and three times less frequently for post-editing REF. This gap was similar in EN-FR, though within the MT systems themselves, the GOOGLE system had a significantly higher error rate for *Mistranslation* errors (38.80 mean) than the next best system, i.e. BASE-NMT (22.73); the contextual MTCUE achieved an even lower rate of 20.10. Interestingly, MTCUE also produced outputs of higher *Fluency* than other systems, even surpassing REF for EN-FR, though insignificantly at the selected confidence interval (80%).

In both language pairs, the *Omission* error was consistently marked the fewest times in GOOGLE-generated text (see **Translation quality** → *Omis-*

sion). In both cases, REF scored significantly above the mean. This is unsurprising: translations authored by the general-purpose GOOGLE engine tend to be overly literal and faithful to the source, while in the domain of dialogue, the HT often needs to let go of individual features of the source text or opt for alternative expressions to maintain the brevity and dynamics of the source dialogue, leading to spontaneous omissions in the reference translations. To exemplify, GOOGLE consistently unnecessarily translated the English “(...), you know,” to “(...), wissen Sie,” in German, necessitating additional post-editing in our study. A similar error was typically avoided by

the other systems, due to their data-learned preference for brevity and dynamically expressive language. As a result, both systems were marked with *Omission* more times than GOOGLE. In fact, MTCUE scored even more *Omissions* than BASE-NMT, suggesting that MTCUE’s omission behaviour more closely matches that of professional HTs. Other **Translation quality** errors were relatively infrequent and with insignificant differences between systems.

To capture context-related issues, we provided categories for the most frequent contextual errors: *Incorrect gender*, *Plural/singular form* and *Formal/informal mismatch*. Since the perception of speaking style in dialogue is subjective and difficult to gauge, we did not provide explicit ways for the PEs to mark speaker style errors to avoid biasing them towards thinking in terms of what is a characteristic way of expression for the given speaker. Instead, we provided loose categories for **Style**, with the intention of collecting measurements of how often the PEs feel the need to alter the style of the translations. Since all of the post-edited content is dialogue, the style of the translation can be directly associated with the style of the speaker’s expression. Our findings regarding some **Context** categories (*Incorrect gender*, *Formal/informal mismatch*) are consistent between the two language pairs, and MTCUE was found to be superior in most categories in both cases, with the overall score for the **Context** category being significant at 80% confidence for EN-FR. The *Plural/singular form* error required few corrections in EN-DE (where BASE-NMT was found superior to MTCUE) and more in EN-FR (where MTCUE was found superior).

The findings from the **Style** category also work in favour of contextual MT, where it was found comparable to non-contextual systems for the EN-DE pair and significantly better than them for the EN-FR pair, requiring the fewest style-based adjustments, even fewer than REF. Within the EN-DE pair, *Subjective style changes* were flagged only up to 4 – 5 times per 100 segments for any system, and a consistent number of times between systems, and *Awkward style* was flagged the fewest times for REF (4.68 on average), much less frequently than for the other systems, among which GOOGLE required the most edits and BASE-NMT the fewest.

Overall, our error count analysis suggests that within the EN-FR pair, MTCUE has significantly

reduced the number of errors marked for contextual and stylistic reasons compared to non-contextual systems, while not degrading overall translation quality. The findings within the EN-DE pair are too variable to yield definitive conclusions but entail no degradation of quality leading from the inclusion of context, a significant improvement for contextual phenomena compared to GOOGLE, and highlight that MTCUE makes the fewest contextual errors overall.

5.2 Analysis of Effort and Quality

This section delves into the analysis of per-PE effort spent post-editing or translating the outputs of each system. Based on the observation that some measurements of editing time and keystrokes were out of the distribution, we normalised these by first computing the 97.5th percentile for the given language pair and task (translation or post-editing) and set all per-segment measurements to be capped at that percentile. Our obtained percentiles were: 37 seconds and 69 keystrokes for translation, and 45 seconds and 54 keystrokes for post-editing.

Effort per PE As per Figure 3, the results for the EN-DE pair suggest that each PE contributed a similar effort. Interestingly, the error rate and effort measures of these PEs are closer in magnitude to the outlier **PE.F3** within the EN-FR pair. Putting PEs from both pairs together we find an interesting correlation: those PEs who expressed a preference for post-editing marked significantly fewer errors overall. We suspect that professionals who expressed a preference for translation opted for spending any effort necessary to match the quality of the resulting text to what they would have produced from scratch, while the post-editing enthusiasts contributed fixed effort, possibly characteristic of their usual post-editing assignments.

The error rate for this pair points to GOOGLE as the system consistently requiring the most edits, and REF the least, though only **PE.G4** made drastically fewer edits to this already production-ready text. Between BASE-NMT and MTCUE, **PE.G2** and **PE.G3** found MTCUE to be less erroneous (and **PE.G3** found it to be on par with REF), while **PE.G1** and **PE.G4** identified fewer errors in BASE-NMT.

According to **PE.G2**, the quality of translations from GOOGLE and BASE-NMT is comparable, requiring the most complex and laborious edits. MTCUE’s hypotheses required less work

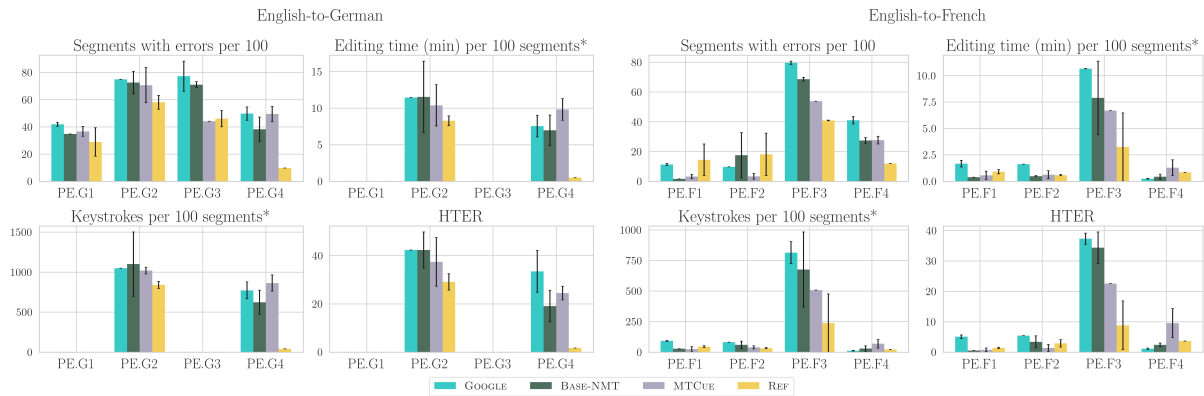


Figure 3: Effort for each PE within both language pairs.

from this PE, and REF text still less. Results obtained from **PE.G4**'s edits are different, revealing next to no edits to the REF text, (which could be interpreted as them being the least subjective of the PEs, only making edits when they are necessary). This PE found MTCUE to require more edits than BASE-NMT and on par with GOOGLE. Interestingly, even though editing MTCUE's outputs took more time and keystrokes, GOOGLE's outputs yielded a HTER value about 10 points higher than MTCUE. Since GOOGLE is the more literal MT system, and MTCUE produces more dialogue-like responses, these findings suggest that, other things being equal, a literal and overly long translation of dialogue may take less effort to post-edit than an incorrect platonic (dialogue-like) response, even if more profound edits are required.

Approach to REF Since the PEs were told about the research nature of the project, they might have approached this project with less vigilance than if the work was undertaken for actual clients. On the flip side, some may have eventually realised they were dealing with some MT outputs – they were not told this explicitly – and became more scrupulous as a result, expecting to make many more corrections than in a typical post-editing task. This would perhaps explain why some PEs took to post-editing REF at rates sometimes matching the outputs of the MT systems, with three of them doing so at a rate of over 40 errors per 100 segments.

Comparison with translation effort In Figure 4 we compare the unnormalised post-editing effort (exclusive of REF) to the FST effort for one episode of the cooking show. For both language pairs, FST required 4 to 6 times the effort of post-editing, by both measures.

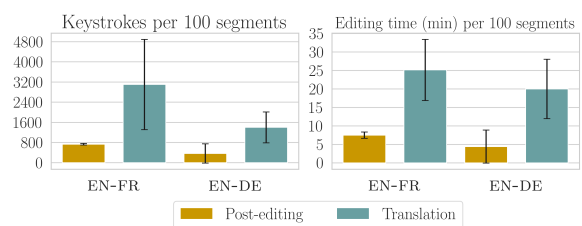


Figure 4: Effort comparison of FST and post-editing MT.

5.3 Analysis of the professionals' views on post-editing and MT

Finally, we present the PEs' responses to a survey regarding views on post-editing and machine translation. Most of the German PEs expressed a preference for FST over post-editing, with three voicing frustration with MT's stiffness and literal nature, omitting aspects of the original text such as slang, gender agreement, references to the video and people's speaking styles. They view translation as a more creative process which can yield idiomatic and fluent translations. They also noted that post-editing currently demands more effort than translating from scratch at times, yet it is compensated at a lower rate than translation. To one PE, post-editing felt like damage control.

Conversely, three out of four French PEs expressed a preference for post-editing, justifying the choice with their specialisation. The fourth PE was dissatisfied with the amount of subtitle formatting errors within our project, commenting that FST would have focused more on content.

PEs in both languages agreed that MT can be a helpful tool, and praised the recent developments, but still concurred that the substantial gap in quality persists, and renders MT insufficiently competent to replace FST. However, they were optimistic about future developments within MT. The

majority of PEs recognized the advantages of post-editing, such as the reduction of temporal effort in some cases and the potential to improve consistency in translating terminology, and enabling greater attention to detail. However, presently these benefits can fail to materialise in practice, emphasising the importance of further work on implementation quality of post-editing workflows.

5.4 Examples of challenges

We present two examples of corrections made in the post-editing process to reflect what kind of corrections required attention as well as what mistakes need to be improved upon in the future.

Example 2		Target: German
Source	No way, no way.	
Video context	<i>The victorious family is in disbelief about their triumph.</i>	
MTCUE (X)	Auf keinen Fall. (‘Under no circumstance.’)	
Post-ed.	Unmöglich. (‘Unbelievable.’)	
Error	Other: inconsistency with video	

Example 2 presents a scenario where MTCUE incorrectly interprets the exclamation *No way* as *Under no circumstance*, which fails to account for the sense of disbelief and amazement that the victorious family is experiencing. Such an interpretation relies strongly on the visual context, of which effective incorporation into the machine translation process in a multi-modal framework is an area for future work.

Example 3		Target: German
Video context	<i>Two cooks and a chopping board.</i>	
Source N	Get that Welly on that board.	
Reference N	Leg das Welly auf das Brett.	
MTCUE (X)	Stell die Welly auf das Brett.	
Post-ed.	Legt das Wellington auf das Brett.	
Error	Awkward style	
Source N+1	She’s on.	
Reference N+1	Es ist drauf.	
MTCUE (X)	Sie ist dran.	
Post-ed.	Ist drauf.	
Error	Other: inconsistency with video	

Example 3 presents a two-error scenario. Firstly, MTCUE uses the incorrect German preposition *an/dran* to translate the English *on*, instead of the correct *auf/drauf* (*on that board* = *auf das Brett*). The more interesting error comes from mis-translating *She* as *Sie*. The pronoun is a reference

to pork Wellington, abbreviated to *Welly* by the speaker, and incorrectly assigned the feminine article *sie*, instead of the neuter *das*. The speaker personifying the pork in Source N+1 (referring to it as *She*) complicates things, and so even a document-level system could have trouble interpreting what *Welly* actually is. The correct interpretation is crucial to selecting the right verb *legen* over *stellen* which should be used to translate *get* when referring to meat. Though it was marked with an *inconsistency with video* error, it is challenging to outline the minimal set of context information sufficient for the correct treatment of this example. The context of cooking, the light-hearted, casual character of the show and the manner of British speech, as well as what meal is being made and what the cooks are doing at the moment, all could aid this process. An important challenge for future contextual systems is going to be to discern which type of information is necessary and when.

6 Conclusions and Future Work

We have presented a case study on post-editing MT of subtitles for TV series in a multi-modal scenario, with a focus on contextual MT. We found that the MT models custom-trained on dialogue required less post-editing effort than the one-size-fits-all Google Translate, potentially due to the overbearing literalness and stiffness of the latter system’s outputs. We also found that some post-editors amended production-approved human translations at high rates, with hypervigilance about dealing with MT as a possible cause. Our results did not determine a significant difference in post-editing effort between MTCUE and BASE-NMT. However, the inclusion of context in MTCUE yielded fewer errors in the **Style**, **Context** and **Fluency** categories, motivating our future exploration of context-inclusive models. We further found that post-editing any MT output required four to six times less technical and temporal effort compared to FST, making it a promising cost-effective venture. However, cognitive effort should be measured in future studies, given the exit survey sentiment that post-editing was sometimes harder and less interesting than FST. Our future experiments will employ larger cohorts of PEs and split them into groups who post-edit non-contextual and contextual inputs exclusively, so that clearer feedback can be collected, as well as to minimise the variance in effort.

7 Acknowledgements

This work was completed as part of Sebastian Vincent's PhD, which was partially funded by ZOO Digital. Sebastian was also supported by the Centre for Doctoral Training in Speech and Language Technologies (SLT) and their Applications funded by UK Research and Innovation (grant number EP/S023062/1).

References

- [Bawden et al.2018] Bawden, Rachel, Rico Sennrich, Alexandra Birch, and Barry Haddow. 2018. Evaluating discourse phenomena in neural machine translation. NAACL HLT 2018 - 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference, 1:1304–1313.
- [C. M. de Sousa et al.2011] C. M. de Sousa, Sheila, Wilker Aziz, and Lucia Specia. 2011. Assessing the post-editing effort for automatic and semi-automatic translations of DVD subtitles. In Proceedings of the International Conference Recent Advances in Natural Language Processing 2011, pages 97–103, Hissar, Bulgaria, September. Association for Computational Linguistics.
- [Freitag et al.2021] Freitag, Markus, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. Transactions of the Association for Computational Linguistics, 9:1460–1474.
- [Gupta et al.2019] Gupta, Prabhakar, Mayank Sharma, Kartik Pitale, and Keshav Kumar. 2019. Problems with automating translation of movie/tv show subtitles. CoRR, abs/1909.05362.
- [Huang and Wang2023] Huang, Jie and Jianhua Wang. 2023. Post-editing machine translated subtitles: examining the effects of non-verbal input on student translators' effort. Perspectives, 31(4):620–640.
- [Karakanta et al.2022] Karakanta, Alina, Luisa Benvogli, Mauro Cettolo, Matteo Negri, and Marco Turchi. 2022. Post-editing in automatic subtitling: A subtitlers' perspective. In Proceedings of the 23rd Annual Conference of the European Association for Machine Translation, pages 261–270, Ghent, Belgium, June. European Association for Machine Translation.
- [Koponen et al.2020] Koponen, Maarit, Umut Sulubacak, Kaisa Vitikainen, and Jörg Tiedemann. 2020. MT for subtitling: User evaluation of post-editing productivity. In Proceedings of the 22nd Annual Conference of the European Association for Machine Translation, pages 115–124, Lisboa, Portugal, November. European Association for Machine Translation.
- [Lison et al.2018] Lison, Pierre, Jörg Tiedemann, and Milen Kouylekov. 2018. OpenSubtitles2018: Statistical rescoring of sentence alignments in large, noisy parallel corpora. In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018), Miyazaki, Japan, May. European Language Resources Association (ELRA).
- [Papineni et al.2002] Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- [Rei et al.2020] Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 2685–2702, Online, November. Association for Computational Linguistics.
- [Sharou and Specia2022] Sharou, Khetam Al and Lucia Specia. 2022. A taxonomy and study of critical errors in machine translation. In Proceedings of the 23rd Annual Conference of the European Association for Machine Translation, pages 171–180, Ghent, Belgium, June. European Association for Machine Translation.
- [Tiedemann and Scherrer2017] Tiedemann, Jörg and Yves Scherrer. 2017. Neural machine translation with extended context. In Proceedings of the Third Workshop on Discourse in Machine Translation, pages 82–92, Copenhagen, Denmark, September. Association for Computational Linguistics.
- [Vincent et al.2022] Vincent, Sebastian T., Loïc Barraud, and Carolina Scarton. 2022. Controlling extra-textual attributes about dialogue participants: A case study of English-to-Polish neural machine translation. In Proceedings of the 23rd Annual Conference of the European Association for Machine Translation, pages 121–130, Ghent, Belgium, June. European Association for Machine Translation.
- [Vincent et al.2023] Vincent, Sebastian, Robert Flynn, and Carolina Scarton. 2023. MTCue: Learning zero-shot control of extra-textual attributes by leveraging unstructured context in neural machine translation. In Findings of the Association for Computational Linguistics: ACL 2023, pages 8210–8226, Toronto, Canada, July. Association for Computational Linguistics.
- [Vincent et al.2024] Vincent, Sebastian, Alice Dowek, Rowanne Sumner, Charlotte Blundell, Emily Preston, Chris Bayliss, Chris Oakley, and Carolina

Scarton. 2024. Reference-less analysis of context specificity in translation with personalised language models. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), Torino, Italia, May. European Language Resources Association (ELRA) and the International Committee on Computational Linguistics (ICCL).

Training an NMT system for legal texts of a low-resource language variety (South Tyrolean German – Italian)

Antoni Oliver and Sergi Álvarez **Egon W. Stemle and Elena Chiocchetti**
Universitat Oberta de Catalunya Eurac Research
{aoliverg, salvarezvid} {egon.stemle, elena.chiocchetti}
@uoc.edu @eurac.edu

Abstract

This paper illustrates the process of training and evaluating NMT systems for a language pair that includes a low-resource language variety. A parallel corpus of legal texts for Italian and South Tyrolean German has been compiled, with South Tyrolean German being the low-resourced language variety. As the size of the compiled corpus is insufficient for the training, we have combined the corpus with several parallel corpora using data weighting at sentence level. We then performed an evaluation of each combination and of two popular commercial systems.

1 Introduction

Neural machine translation (NMT) has shown outstanding performance and translation quality compared to previous models (Bentivogli et al., 2016). However, there is a translation quality gap to fill for low-resource languages (Aranberri and Iñurrieta, 2024; Goyle et al., 2023; Ranathunga et al., 2023) due to the significant amount of parallel data that is required to learn useful mappings between languages (Lakew et al., 2018). Besides, the legal domain poses an increased challenge for NMT (Killman, 2023) given the intricate nature of legal language, the necessity to use precise terminology, and the negative consequences of misunderstanding the legal intent (Quinci and Pontrandolfo, 2023).

Our work deals both with the legal domain and with a low-resource language variety of a

big European language (German), as we translated decrees in the language combination Italian – South Tyrolean German. It shows some of the challenges of NMT in the legal domain, focusing on a low-resource variety. This low-resource variety context is not unique but shared by many language communities in relation to the legal and administrative domain; in Europe, for example, by the German-speaking community in Belgium, the Swedish-speaking community in Finland, the Italian-speaking community in Croatia, the Danish-speaking community in Germany and many more.

After having trained our NMT systems with the LEXB (Contarino, 2021) corpus by Eurac Research and having cleaned and curated the data afterwards, we achieved an increase of, depending on the direction of translation, approximately five to eight full points in BLEU (Post, 2018). This shows the importance of training domain-specific NMT systems with high-quality data, especially for low-resource languages.

1.1 Linguistic situation in South Tyrol

South Tyrolean German is the standard variety of German (Ammon et al., 2016) used in Northern Italy in the Autonomous Province of Bolzano (South Tyrol). German is an officially recognized minority language in South Tyrol. The public administration offices are legally bound to use German next to Italian when dealing with the citizens¹ (Presidential Decree No. 670/1972, Art. 99). All administrative documents, local legislation and material aimed at the general public (e.g. websites of local public institutions) must be available in

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹The small Ladin-speaking minority (20,000 speakers) has also been granted extensive language rights. However, these are generally limited to the valleys of Gherdëina and Badia.

Italian and German.

This officially multilingual institutional regime is implemented by translating texts from German into Italian or vice versa, increasingly using machine translation (De Camillis, 2021). Even though the majority (69%) of the South Tyrolean population is German-speaking (ASTAT – Provincial Statistics Institute, 2021) and today many legal texts are drafted in the minority language, only the Italian version of legal texts is legally binding in case of diverging interpretations (Presidential Decree No. 670/1972, Art. 99). This implies that it is a translated text—either translated or post-edited by a human—that often becomes the legally valid text in South Tyrol.

1.2 Challenges in legal translation

The consequences of mistakes in legal translation may be serious (Mattila, 2018) and include financial loss, legal disputes, infringement of basic human rights (e.g. bad interpreting in a criminal court case). Legal language is therefore considered particularly difficult to translate (Killman, 2023; Mattila, 2018). This is partly due to some specific characteristics of legal language, among others (Gualdo and Telve, 2021; Mattila, 2018):

- specific syntactic features and generally long and complex sentences;
- closeness to general language, with general language words often taking on a specific meaning in the legal context (e.g. ‘trust’);
- terminology that is system-bound and therefore may vary even across legal systems using the same language (e.g. ‘antitrust law’ in the US vs ‘competition law’ in the UK) and may additionally vary in meaning—and translation—also across legal subdomains (e.g. in US banking law, ‘withdrawal’ means the removal of money from a bank but in US criminal law it refers to a person separating themselves from criminal activity²; the first term can be translated with *prelievo*, the second with *dissociazione* in Italian);
- use of abbreviations, acronyms and initialisms;
- formulaic legal phraseology that should not be translated literally.

²<https://thelawdictionary.org/withdrawal/>

All these features are present in the Italian and South Tyrolean German legal languages and pose notable challenges to NMT systems. Chromá (2008) stressed the central role of terminology in legal translation by calculating that between 20% and 29% of legal texts consist of terminology.

1.3 South Tyrolean German and translation

South Tyrolean German has syntactic, grammatical and lexical features that generally characterize it as a Southern German variety and are often shared with the Austrian and/or Swiss standard varieties. Examples are the choice of the auxiliary to form the past tense of some verbs and the use of linking elements within compounds. However, its specific terminology in the domain of law and food as well as a significant influence of Italian clearly distinguishes it from the neighbouring varieties (Ammon et al., 2016). Heiss and Soffritti (2018) and Wiesmann (2019) have shown that terminology is also one of the major machine translation issues in the language combination Italian – South Tyrolean German. A more in-depth error annotation by De Camillis et al. (2023) found that mistranslations and bilingual terminology errors were the most represented error categories when machine-translating South Tyrolean legal texts.

Mistranslations comprise several subcategories of mistakes where the source meaning has been incorrectly transferred to the target language. These include multi-word expressions that have conventional—often non-literal—equivalents like collocations and titles of laws, polysemous words that were disambiguated in the wrong way, occurrences of translations with semantically unrelated words and instances of errors in translating gender-sensitive language. The latter is a known bias of NMT systems (Savoldi et al., 2021). The local South Tyrolean legislation must be inclusive of all genders or at least inclusive of the male and female genders (Provincial Law No. 5/2010, Art. 8). This is achieved by using gender-neutral formulations or terms (e.g. *Lehrperson*, ‘teaching person’) and split forms mentioning both the male and female forms (e.g. *Lehrerinnen und Lehrer*, ‘female and male teachers’) in all language versions. Disrespecting this requirement by generally using only male terms as NMT systems often do (e.g. by translating a gender-neutral expression like *eine Lehrperson* with *un insegnante*, the male form of teacher in Italian) entails a breach of the

law and causes notable post-editing efforts.

Bilingual terminology errors relate to wrongly translated terms in general but also to improper use of terminology pertaining to other legal systems (e.g. *Land* translated with *stato*, ‘state’, rather than *provincia*, ‘province’, because the term refers to a federated state in Germany and Austria but to the Autonomous Province in South Tyrol). The consequences of mistranslating such terminology from German into Italian in South Tyrolean texts are a wrong attribution of competences to the state rather than to the provincial level of governance. Disrespecting the correct terminology does not necessarily make the NMT output impossible to understand. Many South Tyroleans would grasp the meaning of *Tarifvertrag* (‘collective bargaining agreement’ in Germany), even though the correct legal term in South Tyrol is *Kollektivvertrag*. It may also be relatively easy to amend for post-editors with good in-domain knowledge. However, using correct terminology is essential from a legal point of view. Incorrect terms create doubts as to which legal concept is referred to and which legal texts form the legal basis serve as reference. In addition, in a minority language situation, using inconsistent or incorrect legal terminology impairs legal certainty and discriminates against the members of the minority community, as the latter will face additional issues in understanding their legal texts compared to the members of the majority.

2 Experimental part

2.1 The LEXB Italian-German corpus by Eurac Research

Contarino (2021)’s LEXB corpus, a bilingual parallel corpus of Italian and South Tyrolean German, was slightly refined at Eurac Research. It features local and national legislation retrieved from the LexBrowser database³, which gathers laws, decrees, resolutions, collective agreements and other national legal legislation of interest to South Tyrol. The corpus also contains a limited number of bilingual texts not published in the LexBrowser collection, namely 20 national laws and codes (Civil Code, Criminal Code) translated into German, mainly by the provincial Office for Language Issues. This original corpus data has been further cleaned for the current project using MTUOC-clean-parallel-corpus⁴ and rescored with

³<http://lexbrowser.provinz.bz.it/>

⁴<https://github.com/mtuoc/MTUOC-clean-parallel-corpus>

Table 1: Size of the LEXB Italian-German parallel corpus by Eurac Research.

Corpus	Segments	tokens ita	tokens deu
raw	173,530	5,027,663	4,569,333
clean	164,291	4,882,422	4,438,953

Table 2: Size of the Italian-German parallel corpus downloaded from Opus Corpus.

Corpus	Segments
Multiparacrawl	30,337,479
EU rescored	4,936,565

MTUOC-PCorpus-rescorer⁵ (Oliver and Álvarez, 2023). The number of segments and tokens of the raw compiled corpus and the clean and rescored version is included in Table 1. As we can observe, the number of available parallel segments is inadequate for training an NMT system. For this reason, we have combined this corpus with several parallel corpora, as described in the following subsection.

2.2 Other Italian-German corpora used

Table 2 describes the corpora used and their respective sizes in unique segments and tokens. The EU corpus was obtained by concatenating and deduplicating the following corpora: DGT, ELRC-EMEA, EMEA, Europarl and JRC Acquis. The resulting corpus was cleaned and rescored. All the corpora included in this subsection were obtained from Opus Corpus⁶ (Tiedemann, 2009).

2.3 Tools used to train the NMT systems

We used the following tools to train the NMT systems:

- To preprocess the corpora: MTUOC-corpus-preprocessing⁷. This tool allows to use, among other algorithms, sentence-piece⁸ (Kudo and Richardson, 2018).
- Marian NMT⁹ (Junczys-Dowmunt et al., 2018)

⁵<https://github.com/mtuoc/MTUOC-PCorpus-rescorer>

⁶<https://opus.nlpl.eu/>

⁷<https://github.com/mtuoc/MTUOC-corpus-preprocessing>

⁸<https://github.com/google/sentencepiece>

⁹<https://marian-nmt.github.io/>

Table 3: Evaluation results for the Italian-German NMT systems.

System	BLEU	chrF2	TER
	($\mu \pm 95\%$ CI)	($\mu \pm 95\%$ CI)	($\mu \pm 95\%$ CI)
Baseline: Multiparacrawl	37.8 (37.8 \pm 1.6)	58.3 (58.3 \pm 1.0)	57.0 (57.0 \pm 1.7)
EU	29.6 (29.6 \pm 1.2) (p = 0.0010)*	53.4 (53.4 \pm 0.9) (p = 0.0010)*	60.8 (60.8 \pm 1.2) (p = 0.0010)*
EURAC-EU	52.5 (52.5 \pm 2.0) (p = 0.0010)*	68.6 (68.6 \pm 1.1) (p = 0.0010)*	42.0 (41.9 \pm 1.9) (p = 0.0010)*
EURAC-EU-Multiparacrawl	47.9 (47.9 \pm 1.9) (p = 0.0010)*	64.2 (64.2 \pm 1.0) (p = 0.0010)*	45.5 (45.5 \pm 1.5) (p = 0.0010)*
GoogleT	44.1 (44.1 \pm 1.3) (p = 0.0010)*	65.6 (65.6 \pm 0.7) (p = 0.0010)*	45.8 (45.8 \pm 1.2) (p = 0.0010)*
DeepL	36.8 (36.8 \pm 1.2) (p = 0.0849)	63.6 (63.6 \pm 0.7) (p = 0.0010)*	51.3 (51.2 \pm 1.1) (p = 0.0010)*

Table 4: Evaluation results for the German-Italian NMT systems.

System	BLEU	chrF2	TER
	($\mu \pm 95\%$ CI)	($\mu \pm 95\%$ CI)	($\mu \pm 95\%$ CI)
Baseline: Multiparacrawl	33.3 (33.3 \pm 1.1)	56.9 (56.9 \pm 0.9)	54.9 (54.9 \pm 1.1)
EU	47.3 (47.3 \pm 1.4) (p = 0.0010)*	67.9 (67.9 \pm 1.0) (p = 0.0010)*	44.9 (44.9 \pm 1.3) (p = 0.0010)*
EURAC-EU	53.5 (53.5 \pm 1.6) (p = 0.0010)*	71.0 (71.0 \pm 1.0) (p = 0.0010)*	39.1 (39.1 \pm 1.4) (p = 0.0010)*
EURAC-EU-Multiparacrawl	48.3 (48.3 \pm 1.7) (p = 0.0010)*	66.1 (66.1 \pm 1.1) (p = 0.0010)*	44.2 (44.2 \pm 1.5) (p = 0.0010)*
GoogleT	43.5 (43.5 \pm 1.1) (p = 0.0010)*	68.0 (68.0 \pm 0.7) (p = 0.0010)*	44.0 (44.1 \pm 1.1) (p = 0.0010)*
DeepL	47.6 (47.5 \pm 1.3) (p = 0.0010)*	70.0 (70.0 \pm 0.7) (p = 0.0010)*	42.1 (42.1 \pm 1.1) (p = 0.0010)*

2.4 Training procedure

With the corpora described in Sections 2.1 and 2.2 and the tools described in Section 2.3, we trained the following systems in both directions (Italian-German and German-Italian):

- Multiparacrawl: these are the baseline systems trained using the Multiparacrawl corpus (see Section 2.2).
- EU: these systems were trained using the EU corpus (see Section 2.2).
- EURAC-EU: these systems were trained using the EURAC corpus (see Section 2.1) with a sentence weight of 1 and the EU corpus with a sentence weight of 0.5.
- EURAC-EU-Multiparacrawl: these systems were trained using the EURAC corpus with

a sentence weight of 1, the EU corpus with a sentence weight of 0.5 and the Multiparacrawl corpus with a sentence weight of 0.25.

All the corpora have been split into training, validation and evaluation parts. As corpora have been deduplicated, no common segments are present in these subsets. Validation and evaluation sets are formed by 5,000 segments each, and the rest of the segments are used in the training subset. For the EURAC-EU and EURAC-EU-Multiparacrawl corpora, the validation and evaluation subset segments are selected from the EURAC corpus.

All the training processes were performed on a computer with 2 GPUs NVIDIA RTX A 5000 with 24GB each, with the following parameters:

- Guided alignment using eflomal¹⁰ (Östling and Tiedemann, 2016).
- Size of vocabularies: 32,000
- Valid metrics: cross-entropy and bleu-detok
- Patience: 10 on all metrics.
- Type of model: transformer
- Max length of training segments: 150 tokens.

2.5 Evaluation

To evaluate the trained systems, we have used 1,000 segments from the evaluation sets. The trained systems were evaluated along with two popular commercial NMT systems: Google Translate¹¹ and DeepL¹². We accessed both commercial systems through their respective APIs using Python scripts.

For the evaluation, we used three automatic metrics implemented in Sacrebleu¹³ (Post, 2018): BLEU, chrF2 and TER. The appendices present the signatures of the three metrics stating the exact configuration parameters as reported by Sacrebleu.

Tables 3 and 4 show the evaluation results for the Italian-German and German-Italian systems. In both cases, the baseline systems are trained using only the Multiparacrawl corpus. In the evaluation, a paired bootstrap resampling test with 1,000 resampling trials was performed using the `-paired-bs` option in Sacrebleu. In this way, each system is pairwise compared to the baseline system Multiparacrawl. Assuming a significance threshold of 0.05, the null hypothesis can be rejected for p-values < 0.05 (marked with "*" in the tables.)

For both language pairs, the best-performing system according to the used automatic metrics is the systems trained using the EURAC and the EU corpora. For the Italian-German pair, the system improves the baseline system by 14.7 BLEU points, Google Translate by 8.4 BLEU points and DeepL by 15.7 BLEU points. For the rest of the automatic metrics, this system also outperforms the baseline and the commercial systems. The same happens for the German-Italian language pair, where the EURAC-EU system improves the

baseline, Google Translate and DeepL by 20.2, 10 and 5.9 BLEU points, respectively.

Table 5 shows the improvements achieved by the EURAC-EU systems compared with the two commercial systems, Google Translate and DeepL, along with the statistical significance test results for both Italian-German and German-Italian. Figures in the table show the increment of BLEU and chrF2, as well as the decrement of TER, as lower TER values indicate better quality. As we can see in the table, the EURAC-EU systems outperform the commercial systems for the two language pairs and for all the automatic metrics. All these results pass the statistical significance test.

3 Conclusions and future work

Low-resource language situations are challenging for NMT engines. We are working with a low-resource language variety of a major European language and with legal texts, which in itself is a low-resource situation and additionally requires a very particular language.

We have trained an NMT model for the legal domain with the language combination Italian – South Tyrolean German, a low-resource language variety. To this end, we have used and processed a relevant available corpus of legal texts. As this in-domain corpus is not big enough to train NMT systems, we have augmented this data with combinations of other corpora: a corpus created from several EU corpora and Multiparacrawl. The combinations are based on weighting at sentence level, giving higher weight to segments from the compiled in-domain corpus. As a baseline system, we have trained an NMT system using the Multiparacrawl corpus only.

Results show that the best system is the one trained with the in-domain corpus combined with the EU corpora, as it performs better than commercial products for these language combinations. An evaluation was carried out using three of the most frequent assessment metrics (BLEU, chrF2, TER). As positive as these results may seem, a qualitative breakdown of the results, with manual annotations along established criteria (De Camillis et al., 2023) to better understand the specifics of the particular circumstances, is still pending and is planned as the next step.

This paper shows that training tailored NMT systems can be a viable alternative to commercial systems in a low-resource scenario. Even with lim-

¹⁰<https://github.com/robertostling/eflomal>

¹¹<https://translate.google.com/>

¹²<https://www.deepl.com/en/translator>

¹³<https://github.com/mjpost/sacrebleu>

Table 5: Improvements and statistical significance of the EURAC-EU system vs Google Translate and DeepL.

L.P.	System	BLEU	chrF2	TER
		($\mu \pm 95\%$ CI)	($\mu \pm 95\%$ CI)	($\mu \pm 95\%$ CI)
ita-deu	GoogleT	+8.4 (p = 0.0010)*	+3.0 (p = 0.0010)*	-3.8 (p = 0.0010)*
ita-deu	DeepL	+15.7 (p = 0.0010)*	+5.0 (p = 0.0010)*	-9.3 (p = 0.0010)*
deu-ita	GoogleT	+10.0 (p = 0.0010)*	+3.0 (p = 0.0010)*	-4.9 (p = 0.0010)*
deu-ita	DeepL	+5.9 (p = 0.0010)*	+1.0 (p = 0.0010)*	-3.0 (p = 0.0010)*

ited in-domain data, using data from a similar domain and data weighting techniques, the final system can outperform widely used commercial systems. In particular, low-resource varieties of bigger languages tend to be neglected in research and NMT development, even though the consequences of mistranslation may be serious. With its system-bound terminology and phraseology, the legal domain needs particular attention, as it is relevant for legal and translation professionals increasingly using NMT systems and the general public.

Finally, our results emphasize the importance of curated in-domain corpora to align the results of NMT models with those pertaining to situations with more data.

Acknowledgements

This work has been done in the framework of the research and technology transfer agreement between Eurac Research (Bolzano, Italy) and the Universitat Oberta de Catalunya (UOC, Catalonia, Spain).

Appendices:

Metric signatures

- BLEU: nrefs:1 | bs:1000 | seed:12345 | case:mixed | eff:no | tok:13a | smooth:exp|version:2.3.1
- chrF2: nrefs:1 | bs:1000 | seed:12345 | case:mixed | eff:yes | nc:6 | nw:0 | space:no | version:2.3.1
- TER: nrefs:1 | bs:1000 | seed:12345 | case:lc | tok:tercom | norm:no | punct:yes | asian:no | version:2.3.1

References

- Ammon, Ulrich, Hans Bickel, and Alexandra N. Lenz, editors. 2016. *Variantenwörterbuch des Deutschen. Die Standardsprache in Österreich, der Schweiz, Deutschland, Liechtenstein, Luxemburg, Ostbelgien und Südtirol sowie Rumänien, Namibia und Menonitensiedlungen*. de Gruyter, Berlin, 2 edition.
- Aranberri, Nora and Uxoia Iñurrieta. 2024. When minoritized languages encounter MT: perceptions and expectations of the Basque community. *Jostrans – The Journal of Specialised Translation*, (41):179–205.
- ASTAT – Provincial Statistics Institute. 2021. *South Tyrol in Figures*. Provincial Statistics Institute, Bolzano/Bozen.
- Bentivogli, Luisa, Arianna Bisazza, Mauro Cettolo, and Marcello Federico. 2016. Neural versus phrase-based machine translation quality: a case study. In Su, Jian, Kevin Duh, and Xavier Carreras, editors, *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 257–267, Austin, Texas, November. Association for Computational Linguistics.
- Chromá, Marta. 2008. Translating Terminology in Arbitration Discourse. In Bhatia, Vijay K., Christopher N. Candlin, Jan Engberg, and Jane Lung, editors, *Legal Discourse across Cultures and Systems*, pages 309–328. Hong Kong University Press, Hong Kong.
- Contarino, Antonio. 2021. *Neural Machine Translation Adaptation and Automatic Terminology Evaluation: A Case Study on Italian and South Tyrolean German Legal Texts*. Ph.D. thesis, Università di Bologna, Bologna, Italy.
- De Camillis, Flavia, Egon Stemle, Elena Chiocchetti, and Francesco Fernicola. 2023. The MT@BZ corpus: machine translation & legal language. *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 171–180.

- De Camillis, Flavia. 2021. *La traduzione non professionale nelle istituzioni pubbliche dei territori di lingua minoritaria: il caso di studio dell'amministrazione della Provincia autonoma di Bolzano*. Dissertation, Università di Bologna.
- Goyle, Vakul, Parvathy Krishnaswamy, Kannan Girija Ravikumar, Utsa Chattopadhyay, and Kartikay Goyle. 2023. Neural machine Translation for low resource languages.
- Gualdo, Riccardo and Stefano Telve. 2021. *Linguaggi specialistici dell'italiano*. Carocci, Roma.
- Heiss, Christine and Marcello Soffritti. 2018. DeepL Traduttore e didattica della traduzione dall'italiano in tedesco. *inTRAlinea*, 20(1).
- Junczys-Dowmunt, Marcin, Roman Grundkiewicz, Tomasz Dwojak, Hieu Hoang, Kenneth Heafield, Tom Neckermann, Frank Seide, Ulrich Germann, Alham Fikri Aji, Nikolay Bogoychev, André F. T. Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++. In *Proceedings of ACL 2018, System Demonstrations*, pages 116–121, Melbourne, Australia, July. Association for Computational Linguistics.
- Killman, Jeffrey. 2023. Machine translation and legal terminology. Data-driven approaches to contextual accuracy. In Biel, Łucja and Hendrik J. Kockaert, editors, *Handbook of Terminology. Legal Terminology*, volume 3, pages 485–510. Benjamins, Amsterdam / Philadelphia.
- Kudo, Taku and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71.
- Lakew, Surafel M., Marcello Federico, Matteo Negri, and Marco Turchi. 2018. Multilingual Neural Machine Translation for Low-Resource Languages. *IJ-CoL. Italian Journal of Computational Linguistics*, 4(1):11–25, June. Number: 1 Publisher: Accademia University Press.
- Mattila, Heikki E.S. 2018. Legal Language. In Humbley, John, Gerhard Budin, and Christer Laurén, editors, *Languages for Special Purposes: An International Handbook*, pages 113–150. De Gruyter Mouton, Berlin, Boston.
- Oliver, Antoni and Sergi Álvarez. 2023. Filtering and rescoring the CCMatrix corpus for neural machine translation training. In Nurminen, Mary, Judith Brenner, Maarit Koponen, Sirkku Latomaa, Mikhail Mikhailov, Frederike Schierl, Tharindu Ranasinghe, Eva Vanmassenhove, Sergi Alvarez Vidal, Nora Aranberri, Mara Nunziatini, Carla Parra Escartín, Mikel Forcada, Maja Popovic, Carolina Scarton, and Helena Moniz, editors, *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 39–45, Tampere, Finland, June. European Association for Machine Translation.
- Östling, Robert and Jörg Tiedemann. 2016. Efficient word alignment with Markov Chain Monte Carlo. *Prague Bulletin of Mathematical Linguistics*, 106:125–146, October.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels, October. Association for Computational Linguistics.
- Quinci, Carla and Gianluca Pontrandolfo. 2023. Testing neural machine translation against different levels of specialisation: An exploratory investigation across legal genres and languages. *trans-kom*, 16:174–209, July.
- Ranathunga, Surangika, En-Shiun Annie Lee, Marjana Prifti Skenduli, Ravi Shekhar, Mehreen Alam, and Rishemjit Kaur. 2023. Neural Machine Translation for Low-Resource Languages: A Survey.
- Savoldi, Beatrice, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. Gender Bias in Machine Translation. In Roark, Brian and Ani Nenkova, editors, *Transactions of the Association for Computational Linguistics*, volume 9, pages 845–874, Cambridge. Association for Computational Linguistics.
- Tiedemann, Jörg. 2009. News from OPUS: A collection of multilingual parallel corpora with tools and interfaces. In Nicolov, N., K. Bontcheva, G. Angelova, and R. Mitkov, editors, *Recent Advances in Natural Language Processing*, volume V, pages 237–248. John Benjamins.
- Wiesmann, Eva. 2019. Machine translation in the field of law: A study of the translation of Italian legal texts into German. *Comparative Legilinguistics. International Journal for Legal Communication*, 37:117–153.

Implementing Gender-Inclusivity in MT Output using Automatic Post-Editing with LLMs

Mara Nunziatini
Welocalize

mara.nunziatini@welocalize.com

Sara Diego
Welocalize

sara.diego@welocalize.com

Abstract

This paper investigates the effectiveness of combining machine translation (MT) systems and large language models (LLMs) to produce gender-inclusive translations from English to Spanish. The study uses a multi-step approach where a translation is first generated by an MT engine and then reviewed by an LLM. The results suggest that while LLMs, particularly GPT-4, are successful in generating gender-inclusive post-edited translations and show potential in enhancing fluency, they often introduce unnecessary changes and inconsistencies. The findings underscore the continued necessity for human review in the translation process, highlighting the current limitations of AI systems in handling nuanced tasks like gender-inclusive translation. Also, the study highlights that while the combined approach can improve translation fluency, the effectiveness and reliability of the post-edited translations can vary based on the language of the prompts used.

1 Introduction

This paper aims to explore whether LLMs can be effectively utilized for generating gender-inclusive translations. The goal is to determine if this technology can handle the task, or if the expertise of a linguist is still necessary, and to what extent. The challenge lies in the fact that neural machine translation engines frequently fall short in producing gender-inclusive output. When style guides mandate gender-inclusivity in the final translation, post-editors have to make extensive modifications, therefore the MT output is not beneficial for them.

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

We are investigating a multi-step approach to machine translation in which the translation is first produced by an MT engine and is then reviewed by an LLM, to make it gender-inclusive. The goal is for the LLMs to streamline this process and reduce the need for extensive human intervention.

2 The challenges of inclusive writing

Gender-inclusive writing involves using language that does not reinforce traditional gender stereotypes or exclude individuals based on their gender identity. It aims to promote equality and respect for all genders by adopting inclusive terminology and avoiding gendered language whenever possible. Nowadays, gender-inclusive writing is particularly important as societies worldwide strive for greater gender equality and recognition of diverse gender identities.

In this paper, we decided to focus on the translation from English into Spanish for several reasons. Firstly, this language pair poses several gender bias challenges (as we will see later in the paper). Secondly, it is one of the most relevant language pairs from a business perspective for our company. Thirdly, we have highly-trusted internal linguists who are native Spanish speakers and have experience in the translation and post-editing field. Lastly, we have gender-inclusive language style guides available for this language pair, that we used as a starting point for outlining automatic post-editing guidelines. Still, this work is part of an ongoing effort to include additional languages in this experiment.

As mentioned above, gender-inclusive language presents challenges when translating from English into Spanish. This is mainly due to the grammatical structure and inherent gender marking in the Spanish language. Unlike English, where gender-neutral language is more common,

Spanish is a grammatical gender language (Savoldi et al., 2021), that assigns gender to nouns, adjectives, and pronouns. This gender marking extends to articles ('el' for masculine, 'la' for feminine), and even verb conjugations. For example, the English sentence "The doctor saw the patient" can be translated as "El doctor vio al paciente" (masculine doctor, masculine patient), "El doctor vio a la paciente" (masculine doctor, feminine patient), "La doctora vio al paciente" (feminine doctor, masculine patient) or "La doctora vio a la paciente" (feminine doctor, feminine patient).

This inherent gender marking in Spanish makes it challenging to maintain gender neutrality in translations, especially when dealing with professions, titles, and pronouns. Additionally, Spanish has fewer gender-neutral alternatives compared to English, which further complicates the task of creating inclusive translations. Translators are often requested to navigate these linguistic differences while striving to preserve the intended meaning and promote gender inclusivity in the target language.

2.1 Machine Translation and gender-inclusive language

It has been observed that machine translation exacerbates the challenges related to gender-inclusive language since, due to several factors, the raw MT output often contains gender bias (Savoldi et al., 2021). This highlights the need for post-editing and careful consideration of gender-inclusive language.

In our experience, training a machine translation engine to generate gender-inclusive language is challenging due to several reasons:

- MT engines often lack the ability to understand the nuanced context in which gendered language is used, translating based solely on grammar and vocabulary without considering the broader sociocultural implications of gender. This lack of context is often due to the segmentation process that documents go through in order to be translated in Translation Management Systems.
- Different languages have unique grammatical structures and conventions regarding gender. For instance, while English has relatively more gender-neutral options, languages like Spanish assign gender to nouns, adjectives, and

pronouns. This variability makes it difficult to create a one-size-fits-all approach to gender inclusivity in machine translation.

- Different clients have different requirements for gender-inclusive language.
- Machine translation models are trained on large datasets of translated texts. However, these datasets may not always include sufficient examples of gender-inclusive language, leading to biases in the generated translations.
- Using gender-inclusive language often means rephrasing, for example: "gays" → "hombres y mujeres homosexuales". This is especially true if the source text itself includes gender-biased language. Rephrasing requires a deep understanding of context and linguistic subtleties, which can be challenging for machine translation systems.
- We often receive very generic (if any) gender-inclusivity guidelines from clients, which are not detailed enough to train a model.

Overall, training a machine translation engine to generate gender-inclusive language requires addressing these complex linguistic, cultural, and contextual challenges, which may necessitate advanced techniques in natural language processing and extensive fine-tuning of algorithms.

2.2 LLMs for automatic post-editing of gender-biased translations

It appears that LLMs have the potential to be highly effective tools for post-editing tasks. For example, it has been demonstrated that GPT-4 offers promising results on post-editing (Raunak et al., 2023). Besides, LLMs made by large tech companies go through steps which have a goal of minimizing biases in their outputs (Ouyang et al., 2022). We therefore see the identification and fixing of gender-bias issues (whilst translating text) a challenging and very relevant benchmark for judging and comparing LLMs' performance.

Several experiments have been carried out recently to benchmark MT engines and LLMs, and it has been demonstrated that Neural MT engines keep performing better than LLMs (Welocalize, 2023), especially as for accuracy (Vilar et al.,

2023). We think that by using GPT-4 and PaLM2 for automatic post-editing on the raw MT output, we will take advantage of the accuracy delivered by MT engines while improving the translation's fluency with LLMs.

LLMs' ability to understand the context of a text, thanks to being trained on vast and diverse datasets, allows them to make meaningful and contextually appropriate edits. This, combined with their ability to process and edit large volumes of text relatively quickly, makes them a valuable resource for large-scale projects. Also, LLMs can be fine-tuned according to specific guidelines or style guides, including those for gender-inclusive language. We thought that this could make them a potentially valuable tool for enhancing inclusivity in machine translation outputs.

3 Experimental Settings

3.1 Producing the initial translations with MT and LLMs

For this test we utilized content shared by a client which is a globally recognized technology company, and mindful of gender-inclusivity. The content we selected includes text about product integration, technical services, customer support, sales inquiries, cloud solutions, and community interactions. We have chosen this content type as it is written in a way that appeals to all genders, making it an ideal candidate for the test. The language to be used in the translation must be professional, informative, and inclusive, avoiding any gender-biased terms or phrases. This makes it an excellent example of gender-inclusive content in the tech industry. The content was previously translated, therefore we owned the reference human translation.

Firstly, we are interested in producing the initial translations and finding out how the outputs from 5 different systems compare against the human reference translation. This will allow us to choose the best output (output most similar to the reference human translation) to be used as a starting point to generate the gender-inclusive post-edited translation.

For producing the initial translations, we experimented with a subset of 1,000 segments (15,307 words). The systems we used for initial translation generation are:

1. DeepL. We chose this engine since in our experience it is one of the best-performing engines for en>es-ES.

2. GPT-4 (OpenAI, 2023). We chose this system as it has been proved that it consistently performs better than GPT-3.5 (Raunak et al., 2023).
3. PaLM2. State-of-the-art language model that has better multilingual and reasoning capabilities and is more compute-efficient than its predecessor PaLM (Anil et al., 2023). For this exercise we are using text-bison@002 model.
4. DeepL output post-edited by GPT-4.
5. DeepL output post-edited by PaLM2.

In order to enhance the consistency of the assessments and more accurately represent the methodology of general users, our present efforts will concentrate on the zero-shot learning scenario for LLMs, in which the model is not presented with any examples provided by humans. The prompts used to generate the initial translations and the post-edited version of the initial translations can be found in the Appendix (Appendix A and B).

To measure the quality of the 5 outputs, we will compare each one of them against the reference human translation. This will be done by computing COMET (Rei et al., 2020), BLEU (Papineni et al., 2002) and Levenshtein Edit Distance (in our analysis, we normalize this value by the number of characters in the MT output), as these are 3 of the most commonly used reference-based state-of-the-art neural MT quality metrics in the translation industry. The results of this comparison can be found in Paragraph 4.1 (Table 1).

3.2 Performing automatic post-editing to fix gender bias issues in MT output

Secondly, we will perform automatic post-editing with LLMs (GPT-4 and PaLM2), focused solely on fixing gender-bias issues.

In the context of this study, we created a dummy style guide by merging a generic inclusivity writing manual created by our company and a more detailed inclusive writing style guide provided by our client. We therefore asked GPT-4 to transform the resulting style guide into a list of prompts to be used by GPT-4 itself. The list is appended to this paper and was added to the prompt used to perform the automatic post-editing tasks with both LLMs.

We extracted 200 segments from the initial translations, and annotated gender bias issues. We then generated the post-edited gender inclusive translations with GPT-4 and PaLM2. Our internal linguists then evaluated the effectiveness of GPT-4 and PaLM2 in correcting gender bias errors in both Spanish and English texts at segment level, using a labelling system. Labels were “ALL” if all issues were fixed, “PARTIAL” if only some were addressed, and “NONE” if no issues were corrected. The scores can be found in Paragraph 4.2.

4 Results and analysis

4.1 Initial translations with MT and LLMs

Solution	BLEU	PE Distance	COMET
DeepL	49.70	28.00%	0.89
GPT-4	41.47	31.00%	0.88
PaLM2	46.48	31.00%	0.89
DeepL+GPT-4	45.20	30.00%	0.89
DeepL+PaLM2	50.29	28.00%	0.90

Table 1: Quality scores for the initial translations. The “+” sign in the Solution column is to be interpreted as “post-edited by”.

The results in Table 1 suggest that there is no meaningful difference between the 5 different outputs. DeepL+PaLM2 performed the best in terms of translation accuracy and produced an output which is most similar to human reference. However, DeepL alone and the combined approach of DeepL and GPT-4 also performed well.

While GPT-4 and PaLM2 alone performed reasonably well in terms of translation quality, they did not strictly adhere to the prompt we provided for the post-editing step. The internal linguists who carefully reviewed DeepL output post-edited by GPT-4 and DeepL output post-edited by PaLM2 found that both reworked the text more than necessary, to enhance fluency. In many cases, this resulted in the introduction of unnecessary preferential changes, ignoring the part of the prompt stating “Don’t change anything if the Proposed Translation is accurate and fluent”. In fact, these changes didn’t always improve the accuracy or understanding of the text, but rather added a layer of subjective interpretation that was not present in the original text. Moreover, GPT-4 introduced inconsistencies in terminology. For instance, the term “whitepaper”, which was consistently translated by DeepL as “libro blanco”, was sometimes changed by GPT-4 and PaLM2 into different terms such as “documento técnico”, “documento”, “informe blanco”, “informe

técnico” or “documentación técnica”. Other times, it was left unchanged (“libro blanco”) by both LLMs. These inconsistencies can make the job of the post-editor more difficult, as we believe that it is cognitively less demanding and more time-efficient for a reviewer to rectify a recurring terminology inconsistency in a translation than to deal with a single source term translated into the target language in various ways.

In essence, while GPT-4 and PaLM2 showed potential in enhancing fluency, their tendency to introduce unnecessary changes and inconsistencies in terminology raises concerns about their reliability for consistent and accurate translations. Moreover, GPT-4 frequently added the term “Reviewed” at the beginning of the segments, despite the prompt specifically asking for the reviewed text to be returned alone. A similar behavior was already documented in the literature (Zhang et al., 2023) but it came as unexpected since it did not happen in previous tests performed internally by our teams with a similar prompt. This suggests that GPT-4 may have misinterpreted the instructions or overgeneralized from its training data, leading to unnecessary additions to the translated text. This behavior alone unequivocally underscores the continued necessity for human review in the process.

In our process of selecting the most suitable output for our experiment, we chose DeepL’s output. This decision was based on our evaluation of its performance in terms of accuracy, fluency, and consistency of terminology. Furthermore, in a view of adopting this solution in a larger scale scenario, DeepL alone is more cost-effective and time-efficient compared to DeepL reviewed by LLMs. We found that the additional effort required and expense incurred for LLM usage was not justified by a meaningful improvement in quality.

4.2 Automatic post-editing to fix gender-bias issues in MT output

We now use GPT-4 and PaLM2 to review DeepL’s output and make edits solely aimed to ensure that it is gender-inclusive. This means ensuring that the language used does not favor one gender over another and is respectful and inclusive of all genders, following a series of guidelines added to the prompt.

Segment selection and error marking in the initial translation: To ensure an unbiased and

random selection for this experiment, we extracted 200 segments from the initial translation with DeepL. This random extraction ensures a fair and representative sample of the overall text, as it doesn't favor any particular section of the text. The 200 segments were then analyzed by two internal linguists. These individuals are skilled professionals who specialize in language translation and have a keen understanding of gender bias in language.

These internal linguists reviewed each of the 200 segments and marked any gender bias errors, observing the same guidelines that were included in the prompt. These errors could include language that unfairly represents one gender over another, excludes certain genders, or otherwise fails to be inclusive. Out of the 200 segments analyzed, the internal linguists found that 140 of these segments contained one or more gender bias errors. This means that a significant majority of the segments translated by DeepL had issues with gender bias in the translated text. For example, the Spanish equivalent terms for “analyst”, “customer”, “manager”, “developer”, were often used in their masculine form. On the other hand, 60 out of the 200 segments were found to be free of any gender bias errors. This means that these segments were considered by the internal linguists to be gender-inclusive, or simply did not include challenges for gender inclusivity.

Prompting strategy and post-editing by GPT-4 and PaLM2: Both GPT-4 and PaLM2 were then tasked with editing these segments to make them gender-inclusive. This was done using a specific prompt provided in the Appendix of this paper (Appendix C and D), which would have given GPT-4 and PaLM2 guidance on how to approach this task.

The original gender-inclusive language style guides (which we used as a starting point to create the prompt for post-editing) were written in English but included some examples in Spanish. This created a bit of a dilemma when we were trying to decide the language to use for the prompt. Some research had already been done on this topic (Lai et al., 2023; Zhang et al., 2023), and it appears that LLMs perform better with English prompts even if the task and input texts are intended for other languages. Still, we were curious to see if and how the test outcome differs by changing the prompt language. Therefore, we decided to use two different prompts for gender-inclusive post-editing: first the one in English, and then its translation

into Spanish, produced by a professional translator. The reader can find these in Appendix C and D.

Results and discussion: After GPT-4 and PaLM2 had made their edits, the revised segments were given back to the internal linguists for review. The internal linguists then evaluated the changes made by GPT-4 and PaLM2 both with the Spanish and with the English prompt and determined how effectively it had fixed the gender bias errors. The internal linguists used a labelling system to indicate the effectiveness of the LLM's edits (Raunak et al. 2023) at segment level:

- If the LLM had successfully fixed all the gender bias issues in a segment, the internal linguists labelled it as “ALL”.
- If the LLM had only managed to fix some, but not all, of the gender bias issues, the segment would be labelled as “PARTIAL”.
- If the LLM was unable to fix any of the gender bias issues in a segment, the segment was labelled as “NONE”.

This scoring system allowed us to evaluate the effectiveness of using LLMs for post-editing to remove gender bias from machine translations with the English and Spanish prompt.

We noticed that there were two different dimensions that are worth commenting on:

- Quantitative: the number of errors found and fixed by each LLM and
- Qualitative: the quality of the resulting translation.

Quantitative: By looking at Figure 1, we can notice that GPT-4 is more successful than PaLM2 in fixing gender bias issues. GPT-4 was able to identify and fix the majority of gender bias issues. PaLM2 was not as successful, and almost half of the segments with gender bias issues were not fixed or only partially fixed. The above is true both with the English and Spanish prompt.

In fact, the Spanish and English prompt delivered similar results, with the English prompt delivering slightly better results. In more detail, the test results indicated that:

- PaLM2 – the English prompt delivered a slightly better post-edited translation, as the % of segments with gender bias issues that were not fixed at all (“NONE”) is

smaller compared to the post-edited translation delivered with the Spanish prompt.

- GPT-4 – the difference is more meaningful, with almost 80% of the segments with gender-bias issues completely fixed after post-editing with the English prompt, against the 74% with the Spanish prompt.

Based on this data, we can conclude that:

- GPT-4 is more successful than PaLM2 in this task.
- GPT-4 is somewhat more effective at identifying and fixing gender-bias issues when using English prompt compared to Spanish prompt, while changing the prompt language does not make a meaningful difference with PaLM2.
- There is still a clear need for human review as not all segments with gender bias issues were detected and rectified. However, using LLMs helps reducing the number of changes needed.

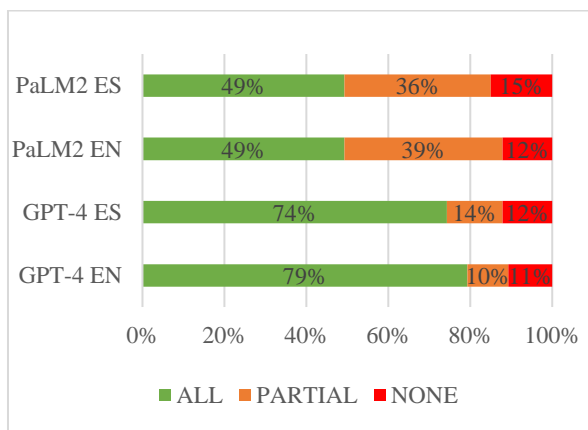


Figure 1: percentage of segments with gender bias errors fixed, partially fixed, and not fixed by PaLM2 and GPT-4, with Spanish (ES) and English (EN) prompts.

Qualitative: The internal linguists noticed that in some cases, the tone of voice was unnecessarily changed in all the four post-edited translations, varying from a formal tone to an informal one. This was against the client’s style guide and was also not requested in the prompt. These unrequested changes can be problematic in a real case scenario, as inconsistencies in tone of voice can complicate the work of the post-editor, who would have to edit much of the text to ensure coherence. Furthermore, in the case of PaLM2, major errors were found with the English prompt, which included neuter forms such as “les desarrolladores”. This solution uses the letter “e” as an alternative for “a” (feminine) or “o” (masculine) in articles,

nouns, and pronouns. It is a recent linguistic development aimed at promoting gender neutrality. However, this solution is not officially recognized (García, 2021b) and, most importantly, it goes against the instructions included in the prompt. Another example is the addition of “sin importar su género” (which translates into “no matter their gender”) in the translation.

Besides, the internal linguists also identified a difference in the quality of the edits between the outputs obtained with the English and Spanish prompts. To address this, we asked our internal linguists to carry out a qualitative ranking of the two translations post-edited by GPT-4 for each segment, judging which gender-inclusive revision was superior from an adequacy and fluency standpoint. We decided to perform this analysis on the translations post-edited by GPT-4 only, without focusing on the translations post-edited by PaLM2, because the former was more successful at this task.

The translator analyzed 140 segments, indicating which between the two post-edited translations demonstrated superior quality for each respective segment. The results indicate that for the greater part of the segments (62%), both translations were comparable from a qualitative standpoint. However, for 22% of the segments, the post-edited translation generated with the English prompt was better, while for the remaining 16%, the post-edited translation generated with the Spanish prompt was better. It was observed that, in those cases where the post-edited translation generated with the English prompt was better, the gender-inclusive solutions proposed were more natural and fluent.

From these results it can be concluded that the choice of the language prompt can have an impact on the quality of the translation, although, in our experiment, in most cases both options delivered similar results.

5 Limitations

The analysis predominantly relies on the outcomes generated by three AI systems, leaving out a comprehensive perspective of the broad array of machine translation systems and large language models available. The study’s focus on a single content type potentially overlooks variations in language use across diverse contents. By examining a limited subset of segments, the study may risk forming a skewed understanding of AI capabilities. Solely focusing on one

language pair fails to consider the inherent structural, complexity, and nuance differences among languages. A more thorough evaluation would require a diverse range of content types and AI systems, a broader selection of segments, as well as multiple language pairs. Finally, we recognize that a thorough comparison of the solutions we examined should ideally include an analysis of output generation speed and associated costs. However, given the page limitations for this paper, we chose to omit this aspect from our current discussion.

6 Conclusions

In conclusion, this study has presented a comprehensive analysis of the performance of LLMs in producing gender-inclusive translations starting from DeepL's raw output. The findings indicate that despite certain potential shown by GPT-4 and PaLM2, the frequent introduction of unnecessary changes, additions, as well as inconsistencies in terminology and tone of voice, raises concerns about their reliability. Furthermore, GPT-4 was found to be more successful than PaLM2 in identifying and fixing gender-bias issues, especially when using an English prompt. This is probably due to the different size of their respective training datasets: GPT-4 was trained on a significantly larger dataset than PaLM2, which means that GPT-4 has "more knowledge" than PaLM2. The study also highlighted the potential impact of the prompt's language on the quality of the translation.

The necessity for human review remains paramount, as not all gender bias issues were detected and rectified by the systems analyzed. Besides, while the use of LLMs to address gender bias issues in translation effectively mitigates the necessity for substantial human intervention in this particular area, it introduces other complications. Specifically, LLMs can create unnecessary alterations in the post-edited translation, such as inconsistencies in terminology and tone of voice. This, in turn, requires further post-editing effort to correct these unintended changes. Therefore, despite the advantages of using LLMs for reducing gender bias, we can't conclusively state that they decrease the overall workload for the post-editor. Further research should delve into the optimization of these systems and their prompts to enhance the accuracy and inclusivity of machine translations.

References

- Anil, R., Dai, A. M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey, P., Chen, Z., Chu, E., Clark, J. H., Shafey, L. E., Huang, Y., Meier-Hellstern, K., Mishra, G., Moreira, E., Omernick, M., Robinson, K., . . . Wu, Y. (2023, May 17). PALM 2 Technical Report. arXiv.org. <https://arxiv.org/abs/2305.10403>
- García, C. (2021b, October 11). El nuevo y tajante mensaje de la RAE sobre el lenguaje inclusivo. La Razón. <https://www.larazon.es/cultura/20211011/flcl3i4owvcwrpvqiqvjljy7wq.html>
- Lai, V. D., Ngo, N. T., Veyseh, A. P. B., Man, H., Deroncourt, F., Bui, T., & Nguyen, T. H. (2023, April 12). ChatGPT Beyond English: Towards a Comprehensive Evaluation of large Language models in Multilingual Learning. arXiv.org. <https://arxiv.org/abs/2304.05613>
- OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Alteschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Zoph, B. (2023, March 15). GPT-4 Technical Report. arXiv.org. <https://arxiv.org/abs/2303.08774>
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike, J., & Lowe, R. (2022, March 4). Training language models to follow instructions with human feedback. arXiv.org. <https://arxiv.org/abs/2203.02155>
- Papineni, Kishore and Roukos, Salim and Ward, Todd and Zhu, Wei-Jing 2002. BLEU: a method for automatic evaluation of machine translation. Proceedings of the 40th annual meeting of the Association for Computational Linguistics, 311-318.
- Raunak, V., Sharaf, A., Wang, Y., Awadalla, H. H., & Menezes, A. (2023). Leveraging GPT-4 for automatic Translation Post-Editing. Findings of the Association for Computational Linguistics: EMNLP 2023. <https://doi.org/10.18653/v1/2023.findings-emnlp.804>
- Rei, Ricardo and Stewart, Craig and Farinha, Ana C and Lavie, Alon. 2020. COMET: A neural framework for MT evaluation arXiv preprint arXiv:2009.09025
- Savoldi, B., Gaido, M., Bentivogli, L., Negri, M., & Turchi, M. (2021). Gender bias in machine translation. Transactions of the Association for Computational Linguistics, 9, 845–874. https://doi.org/10.1162/tacl_a_00401
- Vilar, D., Freitag, M., Cherry, C., Luo, J., Ratnakar, V., & Foster, G. (2022, November 16). Prompting PALM for translation: Assessing strategies and performance. arXiv.org. <https://arxiv.org/abs/2211.09102>
- Welocalize. (2023, August 2). Do LLMs or MT engines perform translation better? - Welocalize. <https://www.welocalize.com/do-llms-or-mt-engines-perform-translation-better/>
- Zhang, B., Haddow, B., & Birch, A. (2023b, January 17). Prompting large language model for machine Translation: A case study. arXiv.org. <https://arxiv.org/abs/2301.07069>

Appendix A. Prompt to generate the Initial Translations

System message: You are a professional translator. You are native English and European Spanish speaker. You specialize in technical translations related to computers, servers, data storage devices, software, and other similar products.

User prompt: Given Source text in English, only return the Translation in European Spanish. Ensure that the translation is fluent and accurately conveys the Source text meaning.

Appendix B. Prompt to post-edit the Initial Translations

System message: You are a professional post-editor. You are native English and European Spanish speaker. You specialize in technical translations related to computers, servers, data storage devices, software, and other similar products.

User prompt: Given Source text in English and its Proposed Translation in Spanish, only return the reviewed translation. Make sure there are no accuracy or fluency issues in the Proposed Translation. If there are, fix them in the reviewed translation. Don't change anything if the Proposed Translation is accurate and fluent.

Appendix C. English prompt to review the Initial Translations and make them gender-inclusive

System message: You are a professional post-editor. You are native English and European Spanish speaker. You specialize in technical translations related to computers, servers, data storage devices, software, and other similar products. You are very interested in inclusive language and always avoid introducing gender bias in your translations.

User prompt: Given the source text in English and its translation into Spanish, only return the post-edited translation. Follow these guidelines:

1. “Check for any gendered terms in the text. If found, can you suggest a gender-neutral alternative for these terms?”
2. “Is the language inclusive for both genders? If not, can you add both gender options, such as ‘bienvenidos/as’ or ‘los/as lectores/as’?”
3. “Can the structure or exact wording of the source text be changed to make the language more inclusive without altering the overall meaning?”
4. “If a gendered term like ‘empleado’ is used, can you think of alternative ways to describe it, such as ‘personal’ or ‘quienes trabajan en...’?”
5. “Is the masculine used as a neutral plural form? If so, can you modify it to avoid sounding awkward?”
6. “Is the ‘pasiva refleja’ used in the text to emphasize the action rather than the subject?”
7. “Are there binary gender representations in the text? If so, can you rewrite it using gender-neutral language?”
8. “Are ‘x’, ‘@’ or ‘e’ used to bypass gender forms? If so, can you suggest an alternative?”
9. “Is a slash (/a) or parentheses (a) used to cover two gender options? If so, can you suggest a different way of doing it?”
10. “Is gender splitting used in the text, i.e., the repetition of masculine and feminine terms? If so, can you suggest a way to avoid it without losing the text’s fluency?”

Note: If none of the above guidelines can be implemented, or when their implementation harms the fluency and naturalness, ask yourself: “Is there a way to maintain the fluency and naturalness of the text while seeking gender neutrality?”

Appendix D. Spanish prompt to review the Initial Translations and make them gender-inclusive

System message: Eres un profesional de la post-edición. Hablas inglés y español de forma bilingüe, y estás especializado en traducciones

técnicas relativas ordenadores, servidores, programas informáticos, y otros productos tecnológicos. Estás muy interesado en el lenguaje inclusivo y siempre evitas introducir sesgos de género en tus traducciones.

User prompt: Dado el texto de origen en inglés y su traducción al español, solo devuelve la traducción post-editada. Sigue estas pautas:

1. “Revisa si hay algún término de género en el texto. Si es así, ¿puedes sugerir una alternativa neutra en género para estos términos?”
2. “¿El lenguaje es inclusivo para ambos géneros? Si no, ¿puedes agregar ambas opciones de género, como ‘bienvenidos/as’ o ‘los/as lectores/as’?”
3. “¿La estructura o el texto exacto del texto fuente pueden ser cambiados para hacer el lenguaje más inclusivo sin alterar el sentido general?”
4. “Si hay un término de género, como ‘empleado’, ¿puedes pensar en formas alternativas de describirlo, como ‘personal’ o ‘quienes trabajan en...’?”
5. “¿Se utiliza el masculino como forma plural neutra? Si es así, ¿puedes modificarlo para que no parezca incómodo?”
6. “¿Se utiliza la ‘pasiva refleja’ en el texto para enfatizar la acción y no el sujeto?”
7. “¿Hay representaciones binarias de género en el texto? Si es así, ¿puedes reescribirlo utilizando un lenguaje neutro en cuanto al género?”
8. “¿Se utilizan ‘x’, ‘@’ o ‘e’ para eludir las formas de género? Si es así, ¿puedes sugerir una alternativa?”
9. “¿Se utiliza una barra (/a) o un paréntesis (a) para cubrir dos opciones de género? Si es así, ¿puedes sugerir una forma diferente de hacerlo?”
10. “¿Se utiliza el desdoblamiento en el texto, es decir, la repetición de términos masculinos y femeninos? Si es así, ¿puedes sugerir una manera de evitarlo sin perder la fluidez del texto?”

Nota: Si ninguna de las pautas anteriores puede implementarse, o cuando su implementación perjudica la fluidez y la naturalidad, pregúntate: “¿Hay una forma de mantener la

fluidez y naturalidad del texto mientras se busca la neutralidad de género?”.

CANTONMT: Cantonese to English NMT Platform with Fine-Tuned Models using Real and Synthetic Back-Translation Data

Kung Yin Hong, Lifeng Han *, Riza Batista-Navarro, Goran Nenadic

Department of Computer Science, The University of Manchester

Oxford Rd, Manchester M13 9PL, United Kingdom

kenrick.kung@gmail.com

{lifeng.han, riza.batista, g.nenadic}@manchester.ac.uk

*corresponding author

Abstract

Neural Machine Translation (NMT) for low-resource languages remains a challenge for many NLP researchers. In this work, we deploy a standard data augmentation methodology by back-translation to a new language translation direction, i.e., Cantonese-to-English. We present the models we fine-tuned using the limited amount of real data and the synthetic data we generated using back-translation by three models: OpusMT, NLLB, and mBART. We carried out automatic evaluation using a range of different metrics including those that are lexical-based (SacreBLEU and hLEPOR) and embedding-based (COMET and BERTscore). Furthermore, we create a user-friendly interface for the models we included in this project, CANTONMT, and make it available to facilitate Cantonese-to-English MT research. Researchers can add more models to this platform via our open-source CANTONMT toolkit, available at <https://github.com/kenrickkung/CantoneseTranslation>.

1 Introduction

Cantonese is one of the most popular dialects of Chinese languages, after the standard language Mandarin (the current official language in China, originally from the Beijing area), originally from the capital of Guangdong province, Guangzhou (a.k.a. Canton) in China. The population of Guangdong province was 129.51 million in 2022 according to the National Bureau of Statistics of China

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹. In addition, Cantonese is also the native language in Hong Kong (HK) and Macau regions which have populations of 7,503,100 and 704,149 in 2023, according to HK Census and Statistics Department² and Macrotrends Global Population statistics.³ Furthermore, because of the economic growth in Guangdong, HK and Macau, many people from other Chinese provinces also learned to speak Cantonese for job purposes and due to cultural influences. There is also a large global population outside of China speaking Cantonese, 85.5 million, according to the Cantonese Language Association (CLA)⁴. In the era of the fast development of natural language processing (NLP), many machine translation (MT) models have been proposed for the majority of languages worldwide. However, *low-resource language MT* remains a challenge for researchers. Cantonese translation using MT, specifically, is under-explored and has not been given much attention thus far.

In this work, we investigate one of the more popular MT methods, i.e. synthetic data augmentation via back-translation and model fine-tuning, as an approach to Cantonese-to-English neural MT (NMT), along the way introducing Cantonese-English as a new language pair. We select several models for evaluation including both smaller and larger language models, and compare their system performance using a range of evaluation metrics. Furthermore, we open-source our toolkit and create a web-based user-friendly platform called **CantonMT** to facilitate research on Cantonese-English translation. A public video demo is available.⁵

¹<https://data.stats.gov.cn/english>

²<https://www.censtatd.gov.hk/en/>

³<https://www.macrotrends.net/global-metrics/countries/MAC/macau>

⁴<https://cantoneseLanguageAssociation.byu.edu/>

⁵CANTONMT demo <https://youtu.be/s8P5fJjS7Ls>

In the next section (Section 2), we survey related work on Cantonese-English MT, data augmentation for MT, and available demos/engines. Section 3 introduces our methodology and framework. Section 4 explains the web-based CANTONMT platform. Section 5 concludes this work with a discussion.

2 Related Work

Research work focussing on Cantonese-English MT has not gained much attention to date. Earliest efforts include the work of (Wu et al., 2006) where example-based and rule-based MT were investigated. In recent years, a project plan on Cantonese-English Translation was put forward by researchers at the University of Hong Kong (HKU) where they proposed to investigate various MT approaches, including rule-based MT (RBMT), example-based MT (EBMT), statistical MT (SMT), gated-recurrent units (GRU) and transformers (Wing, 2020). More loosely related work include research in MT for Cantonese, but without English as the target language. These include dialectal translation between Cantonese and Mandarin Chinese by Zhang (1998), Yi Mak and Lee (Yi Mak and Lee, 2022) and Liu (Liu, 2022).

Data augmentation via backtranslation has been one of the standard practices for generating a synthetic corpus for improving MT performance on low-resource language pairs. This has been popular for both statistical MT (SMT) and NMT (Sugiyama and Yoshinaga, 2019; Graça et al., 2019; Edunov et al., 2020; Nguyen et al., 2021; Pham et al., 2023). However, to the best of our knowledge, none of these efforts focused on Cantonese-to-English translation.

Existing platforms or off-the-shelf demos for Cantonese-to-English MT are very scarce. Popular MT engines from commercial IT companies, including Google Translate⁶ and DeepL Translator,⁷ do not include this language pair. Both of them only included simplified and traditional characters of Mandarin Chinese. Meanwhile, Microsoft Bing Translator⁸ and Baidu, an IT company from China, made the Baidu Translator (Fanyi)⁹ available, which includes Cantonese among several Chinese dialectal languages.¹⁰ In the opposite direction, there

are open-source tools for English-to-Cantonese MT from TransCan.¹¹

3 Experimental Work

We introduce the methodology of CantonMT, experimental evaluations using the initial 38K real bilingual corpus, and extended model evaluations when we acquired 14.5K and 10K more real bilingual data from different sources subsequently.

3.1 Methodology and Framework

The methodology of this work is presented in Figure 1, which includes the following steps:

1. DataPrep: data collection and pre-processing
2. ModelFineTunePhase1: model selection for initial translator fine-tuning (ft, v1)
3. SynDataGenerate: synthetic data generation using the initial translator and cleaned data
4. ModelFineTunePhase2: second step MT fine-tuning using real and synthetic data (ft-syn)
5. ModelEval: model evaluation using both embedding-based metrics (BERTscore and COMET) and lexical metrics (SacreBLEU and hLEPOR)

For data collection, we scraped the data from the public Hong Kong forum LIHKG,¹² which was launched in 2016 and has multiple categories including sports, entertainment, hot topic, gossip, current affairs, etc. We extracted more than 1 million sentences from this website; however, the raw data comes with a lot of noise that needs to be cleaned, an example of which is shown in Figure 5 of Appendix A. We carried out data cleaning to reduce noisy strings as well as data *anonymisation* by removing user IDs from the text. We also filtered out the sentences that were too short, i.e., with less than 10 Chinese characters. In the end, we prepared 200K clean monolingual Cantonese sentences for parallel synthetic data generation purposes. We shuffled the data for model training.

In model fine-tuning phase 1, we aim to train a set of reasonable Cantonese-English MT models for synthetic data generation and model comparisons. The baseline models we selected are OpusMT, NLLB and mBART. These were chosen to

⁶<https://translate.google.com>

⁷<https://www.deepl.com/translator>

⁸<https://www.bing.com/translator>

⁹<https://fanyi.baidu.com/>

¹⁰All these websites were last visited 4th March 2024.

¹¹<https://github.com/ayaka14732/TransCan>

¹²<https://lihkg.com>

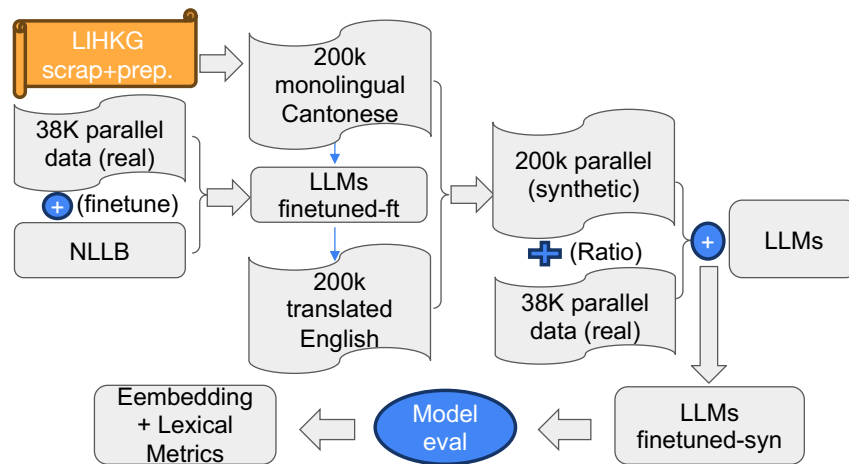


Figure 1: CANTONMT Pipeline: data collection and preprocessing, synthetic data generation, model fine-tuning, model evaluation

answer the following questions: (1) How much does model size impact fine-tuning performance? For this, we use Opus-MT which is a much smaller model trained on the Opus corpus using the MarianMT framework and NLLB-200, a very large language model pre-trained on 200+ languages from Meta-AI; (2) To what extent does it matter if the pre-trained translation models are exposed to Cantonese in their pre-training? For this, we add mBART (mbart-large-50-many-to-many-mmt) which is another LLM but without Cantonese in its pre-training, vs NLLB which includes Cantonese. Because the full-size NLLB is too large, we used the distilled model nllb-200-distilled-600M.

We fine-tuned these models using the available bilingual data from a bilingual Cantonese-English dictionary called “Yue-Dian”,¹³ which is in total 44K in size. We divided this data into training, development and testing sets with 38K, 3K and 3K as their respective sizes, in light of the fact that the shared tasks organised by the Workshop on Statistical Machine Translation (WMT) tend to include around 3K sentences in their test sets.

In Step 3, synthetic data generation, we used the fine-tuned LLMs (LLM-ft-v1) from Step 2 to translate the monolingual Cantonese text we collected in Step 1. In this way, we obtain 200K back-translated English sentences; these synthetic sentences together with the Cantonese sentences create the 200K synthetic parallel corpus we generated. From now on, we will refer to the synthetic parallel corpus as 200K-ParaSyn.

In Step 4, we apply different ratios on the real parallel data we have at hand and on 200K-ParaSyn

for LLM fine-tuning. We also test the influence of model switches, i.e. using different types of LLMs for LLM-ft (Phase 1) and LLM-syn (Phase 2).

In the last step, we deploy the fine-tuned LLMs in Phase 2 (LLM-syn) on the same test data and compare the results with LLM-ft (Phase 1) and baseline models without fine-tuning. We also report comparisons with commercially available translation engines such as the Baidu Translator, Bing Translator and GPT4. The implementation of GPT-4 that we used is Cantonese Companion, which was custom-made for translation to Cantonese by a community builder.¹⁴

We used a range of different evaluation metrics including the lexical-based SacreBLEU (Post, 2018) and hLEPOR (Han et al., 2013a; Han et al., 2021), and the embedding-based BERTscore (Zhang* et al., 2020) and COMET (Rei et al., 2020). hLEPOR has reported much higher correlation scores to the human evaluation than BLEU and other lexical-based metrics on the WMT shared task data (Han et al., 2013b). However, recent WMT metrics task findings have demonstrated the advantages of neural metrics based on embedding space similarities (Freitag et al., 2022).

3.2 Evaluations of CANTONMT

The learning curves of three base models during training using the 38K real data are shown in Figure 2 from left to right for mBART, NLLB-200 and Opus-MT. We used three epochs for mBART because it is too large for the computational resources available to us. From the learning curves, we can

¹³<https://words.hk>

¹⁴<https://chat.openai.com/share/7ee588af-dc48-4406-95f4-0471e1fb70a8>

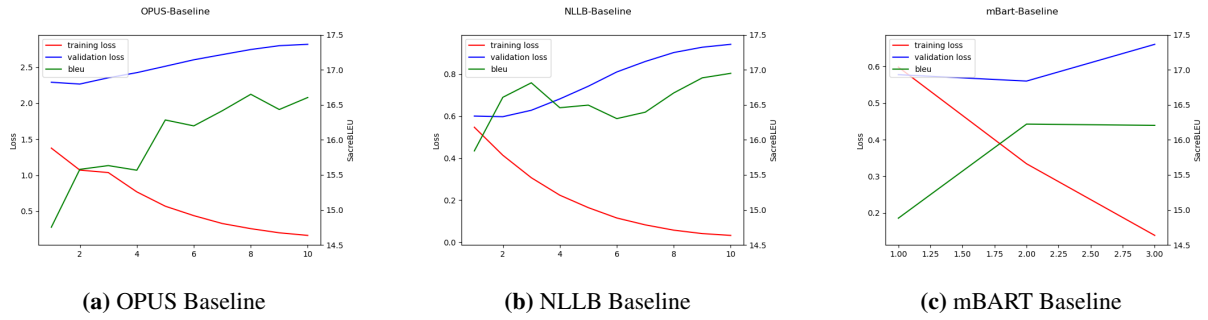


Figure 2: Learning curves during model training using real data.

Model Name	SacreBLEU	hLEPOR	BERTscore	COMET
nllb-forward-bl	16.5117	0.5651	0.9248	0.7376
nllb-forward-syn-h:h	15.7751	0.5616	0.9235	0.7342
nllb-forward-syn-1:1	16.5901	0.5686	0.925	0.7409
nllb-forward-syn-1:1-10E	16.5203	0.5689	0.9247	0.738
nllb-forward-syn-1:3	15.9175	0.5626	0.924	0.7376
nllb-forward-syn-1:5	15.8074	0.562	0.9237	0.7386
nllb-forward-syn-1:1-mbart	16.8077	0.571	0.9256	0.7425
nllb-forward-syn-1:3-mbart	15.8621	0.5617	0.9246	0.7384
nllb-forward-syn-1:1-opus	16.5537	0.5704	0.9254	0.7416
nllb-forward-syn-1:3-opus	15.9348	0.5651	0.9242	0.7374
mbart-forward-bl	15.7513	0.5623	0.9227	0.7314
mbart-forward-syn-1:1-nllb	16.0358	0.5681	0.9241	0.738
mbart-forward-syn-1:3-nllb	15.326	0.5584	0.9225	0.7319
opus-forward-bl-10E	15.0602	0.5581	0.9219	0.7193
opus-forward-syn-1:1-10E-nllb	13.0623	0.5409	0.9164	0.6897
opus-forward-syn-1:3-10E-nllb	13.3666	0.5442	0.9167	0.6957
baidu	16.5669	0.5654	0.9243	0.7401
bing	17.1098	0.5735	0.9258	0.7474
gpt4-ft(CantoneseCompanion)	19.1622	0.5917	0.936	0.805
nllb-forward-bl-plus-wenlin14.5k	<i>16.6662</i>	<i>0.5828</i>	<i>0.926</i>	<i>0.7496</i>
mbart-forward-bl-plus-wenlin14.5k	15.2404	0.5734	0.9238	0.7411
opus-forward-bl-plus-wenlin14.5k	13.0172	0.5473	0.9157	0.6882
nllb-200-deploy-no-finetune	11.1827	0.4925	0.9129	0.6863
opus-deploy-no-finetune	10.4035	0.4773	0.9082	0.6584
mbart-deploy-no-finetune	8.3157	0.4387	0.9005	0.6273
nllb-forward-all3corpus	<i>16.9986</i>	<i>0.583</i>	<i>0.927</i>	<i>0.7549</i>
nllb-forward-all3corpus-10E	16.1749	0.5728	0.9254	0.7508
mbart-forward-all3corpus	16.3204	0.5766	0.9253	0.7482
opus-forward-all3corpus-10E	14.4699	0.5621	0.9191	0.7074

Table 1: Automatic Evaluation Scores from Different Models in CANTONMT. bl: bilingual real data; syn: synthetic data; h:h - half and half; 1:1/3/5 - 100% real + 100/300/500% synthetic; 10E: 10 epochs (default: 3); top-down second slot: model switch: model type using NLLB but synthetic data from other models (mBART and OpusMT); top-down third slot: including model switch for mBART fine-tuning using synthetic data generated from NLLB; similarly top-down forth slot: including model switch for OpusMT fine-tuning using synthetic data from NLLB. Bottom slot of Cluster 1: Bing/Baidu Translator and GPT4-finetuned Cantonese Companion; **bold** case is the best score of the same slot among the same model categories. Cluster 2: bilingual fine-tuned models using 38K words.hk data plus 14.5k Wenlin data; *italic* indicates the number outperforms the same model fine-tuned with less data 38K. Cluster 3: Deployed Model without fine-tuning Cluster 4: Finetuned with the previous 2 corpora and an additional 10K data from OPUS Corpora we managed to find in the end - it shows the evaluation improvement continues.

see that NLLB-200 has a peak score at epoch 3 then there is a dramatic drop until epoch 6, followed by an increase until epoch 10. In contrast, the Opus-MT model achieves a steady increase in its SacreBLEU score with more epochs, although there are little drops in between.

The automatic evaluation scores from CANTONMT models and other commercial engines are listed in Table 3. Below are some interesting findings from the evaluation outcomes.

- LLM-ft vs -LLM-ft-syn: (1) NLLB-syn-1:1 has slightly better scores than NLLB-bl on all metrics, but increasing the ratio of synthetic data will decrease the scores such as in the 1:3 and 1:5 configurations, with around 1 absolute SacreBLEU point. (2) Similarly, mBART-syn-1:1 also outperforms mBART-ft but increasing the ratio of synthetic data will reduce the evaluation scores such as in the 1:3 configuration. (3) Surprisingly, the synthetic model for Opus-mt does not outperform Opus-ft-bl, which indicates that the quality of the generated synthetic data matters.
- Model Switching Matters: (1) the NLLB fine-tuned model using synthetic data from mBART (second model from the top of the table) produced higher scores than using the synthetic data generated from its own (first model from the top of the table). (2) mBART fine-tuned using NLLB-generated synthetic data also outperforms mBART fine-tuning using only bilingual real data. (3) In a similar situation, Opus-MT performs differently in comparison to the other two models.
- Commercial MT models: (1) GPT4-finetuned produced the highest evaluation scores but the free version of GPTs restricts the input number of strings; the data size used for fine-tuning GPT-4 is unknown and such data is not publicly available to researchers; furthermore, it is unclear how GPT-4 performs MT; in addition, there are risks to data privacy when users choose to use engines from commercial companies. In contrast, CANTONMT is open-source, free, and researchers can continue to fine-tune it with their data or include more models, and is fully *confidential* for users. 2) Bing and Baidu translators produced similar evaluation scores to the best system from CANTONMT, though Bing produced slightly

higher scores than Baidu, especially on the lexical-based metrics SacreBLEU and hLEPOR.

- Comparing to Model Deployment without Finetuning: in Cluster 3 (bottom) of Table 3, model deployment without fine-tuning has much lower scores; these scores show that fine-tuning and synthetic data augmentation lead to a large increase in scores of around 50% for all models using SacreBLEU.

3.3 Adding More Real Data

In the extension of our work, we managed to fine-tune the baseline models using more real data from another source called Wenlin¹⁵ where we obtained another 14.5K parallel Cantonese-English dictionary. We are curious about the model performance using more real data in addition to the 38K training corpus from words.hk. We listed the comparison scores in the second cluster of Table 3 where it shows that the newly fine-tuned NLLB-200 using 52.5K data (38+14.5K=52.5K) produced higher scores on all metrics in comparison to 38K trained model; mBART fine-tuned using 52.5K obtains better scores on three metrics except for SacreBLEU; Opus-MT surprisingly did not get any increase across the metrics. Nevertheless, these outcomes demonstrated the possibility of improving model performance with more available real data, at least for the NLLB and mBART models. Moreover, **data quality matters**: simply adding 14.5K real data to fine-tune NLLB produced higher scores (underlined scores) than the best synthetic system that used 38x2=76K data. Subsequently, when we managed to get another 10K real data from Opus corpus, it shows continuous improvement by training using all three corpus we have, located in the Cluster 4 bottom of the table.

4 CANTONMT Platform

To further facilitate Cantonese-English MT research and for users to easily access freely available fine-tuned models, we developed a user-friendly interface for the CantonMT platform. Users can choose different models and translation directions (Cantonese \Leftrightarrow English) via the interface (Figure 4 in the Appendix). The web application contains two main parts, the Interface and the Server.

¹⁵<https://wenlin.com/>

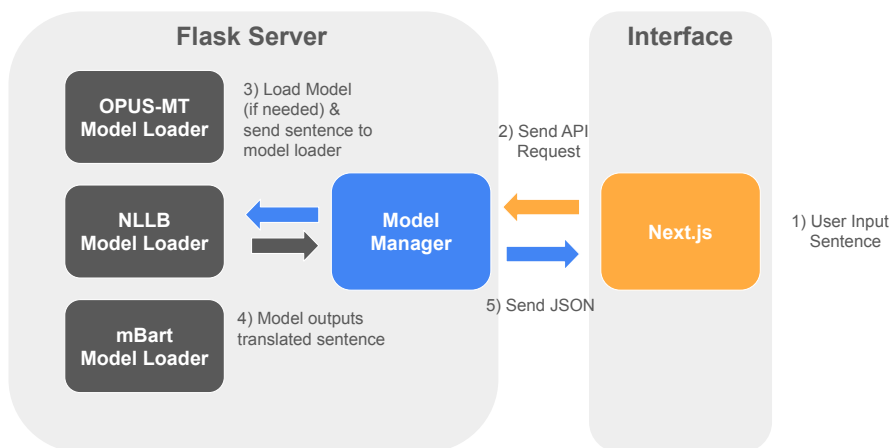


Figure 3: CANTONMT Server and Interface Flowchart diagram.

4.1 User Interface

To test the user interface and different models for translation, users can choose from different model types and source languages, which dynamically capture the available models in the server, and allow users to select different training methods for the model. One can then type the source sentence in the input box and click the “Translate” button to obtain the translation output from the model. The application layout is quite modular in case different model types or languages are added to the system, which could potentially be used as a base framework for different translation systems. It is possible to simply add more languages to the input and output if one wishes to expand the implementations. The look-and-feel of this web application is based on a template (Wrigley, 2023) for an AI Code translator, which was customised and developed in TypeScript with the Next.js framework. The reason for choosing this framework is that it provides a very modern and minimalistic approach.

4.2 Server

A diagram outlining the modules can be seen in Figure 3 to understand the general structure of the server. Users can easily run the server on their local machines by following the instructions provided in a README file. The server has two main functionalities, where the first one will output the list of model paths given the model type and source languages. The second one provides the translation, where one could provide the details of the model and also the sentence in the language specified, and the server would respond with the translated sentence using the model output.

During our implementation, due to memory con-

straints, the server crashed multiple times on our local machine. To mitigate the risk of server crashes, a *model manager* was produced, which implements a Least Recently Used (LRU) cache for the different model loaders, where the least recently used model will be deleted from memory if it exceeds the limit of the number of models. The server is built entirely based upon the *Python Flask* library. The reason for choosing this framework is that the models can be run on top of the Python Transformers library, which provides seamless implementation without much additional effort.

5 Discussion and Conclusion

In this work, we investigated the back-translation methodology for bilingual synthetic data generation for the sake of data augmentation for NMT, on a new language translation direction, Cantonese-to-English. We tested both smaller-sized OpusMT and extra-large LLMs NLLB and mBART both using available bilingual real data and larger synthetic data. Our experiments show that all the fine-tuned models outperformed the baseline deployment models with large margins. Furthermore, the synthetic model nllb-syn-1:1-mbart produced higher scores using the model switch method compared to those without the model switch. Lastly, the best performing fine-tuned models have similar (or even higher) evaluation scores than the current commercially available translators of Baidu and Microsoft-Bing.

In terms of concerns of **data privacy** such as handling of sensitive data (e.g., in clinical applications related to health analytics of patient data (Han et al., 2024; Han et al., 2022)), CANTONMT can be fully controlled by users without interference from any third parties. We open-source our platform so that

researchers can continue to integrate new models into the toolkit to promote Cantonese-English MT. We also plan to carry out human evaluations on the outputs from different systems to get more insights into the system errors.

Limitations

The synthetic data generated in this work is based on the fine-tuned model using 38K words.hk bilingual dictionary corpus, the first corpus we managed to find. This restricted the synthetic data quality. In the following-up work, we plan to use the further fine-tuned model on all three corpora, words.hk, wenlin, and opus-10K, to generate better back-translated synthetic data. We expect this will improve the synthetic data fine-tuned models.

The whole procedure of how difficult it was to collect real Cantonese-English bilingual data shows that Cantonese-English MT is still at its beginning stage with many obstacles and challenges to public research.

Acknowledgements

We thank words.hk and wenlin.com for the data. LH and GN thank the support from the following grants: (1) “Assembling the Data Jigsaw: Powering Robust Research on the Causes, Determinants and Outcomes of MSK Disease”. This project has been funded by the Nuffield Foundation (www.nuffieldfoundation.org), but the views expressed are those of the authors and not necessarily of the Foundation; and (2) the UKRI/EPSC Grant EP/V047949/1 “Integrating hospital outpatient letters into the healthcare data space”.

References

- [Edunov et al.2020] Edunov, Sergey, Myle Ott, Marc’ Aurelio Ranzato, and Michael Auli. 2020. On the evaluation of machine translation systems trained with back-translation.
- [Freitag et al.2022] Freitag, Markus, Ricardo Rei, Nittika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- [Graça et al.2019] Graça, Miguel, Yunsu Kim, Julian Schamper, Shahram Khadivi, and Hermann Ney. 2019. Generalizing back-translation in neural machine translation.
- [Han et al.2013a] Han, Aaron Li-Feng, Derek F. Wong, Lidia S. Chao, Liangye He, Yi Lu, Junwen Xing, and Xiaodong Zeng. 2013a. Language-independent model for machine translation evaluation with reinforced factors. In *Proceedings of Machine Translation Summit XIV: Posters*, Nice, France, September 2-6.
- [Han et al.2013b] Han, Aaron Li-Feng, Derek F. Wong, Lidia S. Chao, Yi Lu, Liangye He, Yiming Wang, and Jiaji Zhou. 2013b. A description of tunable machine translation evaluation systems in WMT13 metrics task. In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 414–421, Sofia, Bulgaria, August. Association for Computational Linguistics.
- [Han et al.2021] Han, Lifeng, Irina Sorokina, Gleb Erofeev, and Serge Gladkoff. 2021. cushLEPOR: customising hLEPOR metric using optuna for higher agreement with human judgments or pre-trained language model LaBSE. In *Proceedings of the Sixth Conference on Machine Translation*, pages 1014–1023, Online, November. Association for Computational Linguistics.
- [Han et al.2022] Han, Lifeng, Gleb Erofeev, Irina Sorokina, Serge Gladkoff, and Goran Nenadic. 2022. Examining large pre-trained language models for machine translation: What you don’t know about it. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 908–919, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- [Han et al.2024] Han, Lifeng, Serge Gladkoff, Gleb Erofeev, Irina Sorokina, Betty Galiano, and Goran Nenadic. 2024. Neural machine translation of clinical text: An empirical investigation into multilingual pre-trained language models and transfer-learning. *Frontiers in Digital Health*, 6:1211564.
- [Liu2022] Liu, Evelyn Kai-Yan. 2022. Low-resource neural machine translation: A case study of Cantonese. In *Proceedings of the Ninth Workshop on NLP for Similar Languages, Varieties and Dialects*, pages 28–40, Gyeongju, Republic of Korea, October. Association for Computational Linguistics.
- [Nguyen et al.2021] Nguyen, Xuan-Phi, Shafiq Joty, Thanh-Tung Nguyen, Kui Wu, and Ai Ti Aw. 2021. Cross-model back-translated distillation for unsupervised machine translation. In *International Conference on Machine Learning*, pages 8073–8083. PMLR.
- [Pham et al.2023] Pham, Nghia Luan, Van Vinh Nguyen, and Thang Viet Pham. 2023. A data augmentation method for english-vietnamese neural machine translation. *IEEE Access*, 11:28034–28044.

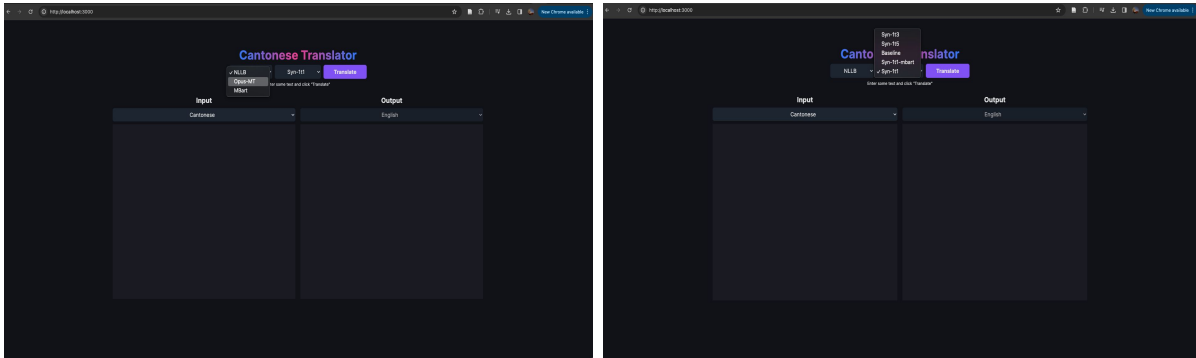


Figure 4: CANTONMT Platform with options of model types, training categories, and translating directions. Frontend: TypeScript with Next.js. Backend: Python - Flask

```
number,date,uid,probation,text,upvote,downvote,postid,title,board,collection_time
#386,2023年11月21日 09:41:17,/profile/[profile-id],FALSE,電視劇得唔得?Game of thrones
red wedding,,3558451,不劇透：邊套戲你睇過有最強twist位???,影視台,
2023-11-21T10:29:59.718892Z
#1,2023年11月20日 14:19:14,/profile/[profile-id],FALSE,"發展商積極推售新盤搶佔市場購買力，永
義集團旗下何文田窩打老道已屆現樓的「譽林」(13日)落實首輪銷售安排，將於(17日)以先到先得形式，發售首張價
單全數30伙。扣除家具優惠及最高折扣後，折實售價由529.7萬元起，折實平均呎價20,935元。
「譽林」上周五發售的30伙，實用面積介乎260至754方呎，戶型涵蓋開放式至三房。價單定價由598.3萬至
1,913.7萬元，呎價介乎20,300元至25,450元。扣除家具折扣優惠及最高樓價10%折扣後，單位折實售價由529.7
萬至1,701.9萬元，折實呎價介乎17,909元至22,612元。
最後結果：",,,3558452,何文田譽林上周五首輪開售30伙 成功售出4伙,房屋台,
2023-11-21T10:30:07.323742Z
```

Figure 5: Example text extracted from LIHKG website with lots noise before cleaning and anonymisation

[Post2018] Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October. Association for Computational Linguistics.

[Rei et al.2020] Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.

[Sugiyama and Yoshinaga2019] Sugiyama, Amane and Naoki Yoshinaga. 2019. Data augmentation using back-translation for context-aware neural machine translation. In *Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019)*, pages 35–44, Hong Kong, China, November. Association for Computational Linguistics.

[Wing2020] Wing, Liu Hey. 2020. Machine translation models for cantonese-english translation project plan.

[Wrigley2023] Wrigley, Mckay. 2023. ai-code-translator.

[Wu et al.2006] Wu, Yan, Xiukun Li, and Caesar Lun. 2006. A structural-based approach to Cantonese-English machine translation. In *International Journal of Computational Linguistics & Chinese Language Processing, Volume 11, Number 2, June 2006*, pages 137–158, June.

[Yi Mak and Lee2022] Yi Mak, Hei and Tan Lee. 2022. Low-resource nmt: A case study on the written and spoken languages in hong kong. In *Proceedings of the 2021 5th International Conference on Natural Language Processing and Information Retrieval, NLPPIR '21*, page 81–87, New York, NY, USA. Association for Computing Machinery.

[Zhang* et al.2020] Zhang*, Tianyi, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

[Zhang1998] Zhang, Xiaoheng. 1998. Dialect MT: A case study between Cantonese and Mandarin. In *COLING 1998 Volume 2: The 17th International Conference on Computational Linguistics*.

A Appendix

Example raw text extracted from LIHKG website can be seen in Figure 5 before cleaning. CantonMT

user-friendly interface is shown in Figure 4 (Frontend: TypeScript with Next.js. Backend: Python - Flask).

The model parameters from OpusMT, extra-large NLLB and mBART are shown in Table 2, which shows that NLLB and mBART have doubled the number of transformer layers and have almost 10 times more parameters than OpusMT.

	Opus	NLLB	mBart
Layers	12	24	24
Hidden Unit	512	1024	1024
Model Parameters	77.9M	615M	610.9M
Language Pair	No	Yes	No
Release Year	2020	2022	2020

Table 2: Parameters from deployed models. Language pair: if the model contains Cantonese-English as a language pair

Explanation of Abbreviations used in the scoring table:

- “nllb-forward-bl”: NLLB fine-tuned model in the forward translation direction (Cantonese-English) using the real 38K bilingual corpus
- “nllb-forward-syn-h:h”: NLLB fine-tuned model using forward-translation generated synthetic data to substitute half of the 38K real data, i.e. 19K real and 19K synthetic
- “nllb-forward-syn-1:1”: NLLB fine-tuned model using forward-translation generated synthetic data with the ratio 1:1, i.e. 38K real and 38K synthetic
- “nllb-forward-syn-1:1-10E”: the same with above corpus setting but running 10 epochs, default is 3 epochs only
- “nllb-forward-syn-1:1-mbart”: NLLB model fine-tuning using forward-translation generated synthetic data by another model mBART, 38K real and 38K synthetic

Model Name	SacreBLEU	BERTscore	COMET
mBART-back-bl	20.3841	0.7944	0.8095
mBART-back-syn-1:1-NLLB+	20.1923	0.7921	0.8068
nllb-back-bl	18.4713	0.7877	0.7927
nllb-back-syn-1:1	17.9400	0.7807	0.7772
nllb-back-syn-1:3	12.0352	0.7628	0.7493
opus-back-bl	18.1496	0.7811	0.7816
opus-back-syn-1:1-NLLB+	17.9346	0.7781	0.7715

Table 3: Automatic Evaluation Scores from Different Models in CANTONMT.

Hyperparameter	Value
Learning Rate	10_{-4}
Weight Decay	0.01
FP16	True

Table 4: Fine-tuning Hyperparameters w/ Hugging Face Trainer API

Advancing Digital Language Equality in Europe: A Market Study and Open-Source Solutions for Multilingual Websites

Andrejs Vasiljevs

Tilde, Vienības gatve 75a, Riga, Latvia
University of Latvia, Raiņa bulv. 29, Riga, Latvia
andrejs@tilde.lv

Neli Vacheva

IDC Bulgaria, Sofia 1040, Bulgaria 36, Dragān Tzankov blvd.
nvacheva@idc.com

Rinalds Viksna

Tilde, Vienības gatve 75a, Riga, Latvia
University of Latvia, Raiņa bulv. 29, Riga, Latvia
rinalds.viksna@tilde.lv

Andis Lagzdīņš

Tilde, Vienības gatve 75a, Riga, Latvia
andis.lagzdins@tilde.lv

Abstract

The paper presents findings from a comprehensive market study commissioned by the European Commission, aimed at analysing multilinguality of European websites and automated website translation services across various sectors. The findings show that the majority of websites offer content in one or two languages, while only less than 25% of European websites provide content in 3 or more languages. Additionally, we introduce Web-T, a collection of open-source solutions facilitating automated website translation with a help of free MT service eTranslation provided by the European Commission and possibility to integrate other MT providers. Web-T solutions include local plug-ins for Content Management Systems, universal plug-ins, and an MT API Integrator, thus contributing to the broader goal of digital language equality in Europe.

1 Introduction

Within the European Union, a diverse linguistic landscape is comprised of 24 official languages and more than 60 regional and minority languages. Several research studies (Pastor et al., 2017; Rehm et al., 2020; Rehm and Way, 2023) and official resolutions (European Parliament, 2018; European Commission, 2008) have underscored a stark discrepancy in the

technological support available for Europe's multitude of languages.

This is especially pertinent in light of the current lack of multilinguality on many European websites (a website of a company, based in Europe, regardless of whether they belong to a European subsidiary of a global corporation or are headquartered in Europe), highlighting the need to promote the use of language technologies to make digital content and online services multilingual and more accessible for all European citizens.

The challenge of limited multilingual support on websites extends beyond Europe and has been highlighted by research in other parts of the world (Wright, 2004; Miraz et al., 2013; Singh et al., 2016; Sargent and Lommel, 2019; Kelly-Holmes, 2019). Supporting a website in multiple languages (translating both UI and the content) can be a time-consuming and expensive process. Automated translations have revolutionised website localization, making it more accessible to businesses, including smaller enterprises and individuals. Despite limitations of automated translations that may not always accurately convey the intended message or account for cultural differences, businesses can benefit from cost savings, speed, and scalability, which allow them to expand their global presence.

In many cases, a precise translation is needed (government, legislation, healthcare, industry specifications, brand identity, etc.), and to reach that high quality of translation, so far, the automated translation must be followed by post-editing by humans. However, automated translation is enough in cases when the users need to get a

general understanding of the content of a webpage and non-perfect translation won't have critical consequences. The automated translation accelerates the translation process and as far as the language models can be trained on specific topics, languages, jargons, and dialects, the quality of the translations can be improved to a level that requires minimum or no human intervention, making the process even more productive (Stasimioti et al., 2020).

Although there are numerous automated translation solutions provided by market players, their use is not dominant on the European web space. To assess and improve the situation, the European Commission (EC) has commissioned an extensive market study on multilingualism of websites in Europe and the development of solutions to support the use of automated translations on websites. The project is implemented in the scope of the Digital Europe Programme's Strategic Objective 5, "Accelerating best use of technologies," and aims to enhance language technologies' capacity within the European public sector and their broader deployment across public and private sectors, NGOs, and academia (European Commission, 2021).

In this paper we present the project findings in analysing the language diversity on the European web space, the use of solutions ensuring automated website translation, and the machine translation services underpinning these solutions.

We also introduce a collection of open-source solutions developed under the project, collectively known as Web-T. These solutions offer free-of-charge automated website translations utilizing the European Commission's eTranslation machine translation (MT) service¹ and are adaptable for integration with other MT providers.

2 Assessing Multilingualism of European Websites

According to a recent study by IDC, there are slightly more than 1 million websites managed by public sector enterprises in Europe (EU 27 plus Albania, North Macedonia, Switzerland, Serbia, Montenegro, and Bosnia and Herzegovina) and about 8.4 million websites managed by private sector entities in Europe.

To assess the multilinguality of European websites, two randomised sample lists of European

websites were compiled. The lists of websites for analysis were compiled by combining a list of websites per country available on builtwith.com, lists of small and medium enterprises (SMEs) provided by national registries, trade organisations, and lists of government institutions, universities, schools, and healthcare institutions on a national level and a regional/city level for each country. A random subset of the lists was used for analysis. One list contains websites sampled from domains of the largest economies: Germany, France, Italy, Austria, and the Netherlands. The list is balanced to include about 20% public sector and about 55% SMEs, with the rest being large or medium companies. The second list contains links to websites of enterprises in EU 27 member states and is balanced between big companies, SMEs, and the public sector.

The Multilingualism Scoring Tool (Viksna et al., 2022) was used to measure the multilingualism of a website. It analyses the textual content of the website and identifies the number of languages used, the distribution of content in various languages, and the presence of multilingual features. Multilingual features are website features that point to this webpage being available in other languages and offering user access to this content, such as language switcher tool/button/link, machine-readable links to translated content, or blocks of text in various languages available for display using JavaScript.

From the first list of largest economies, 426 websites were crawled with a depth of 2 links. Most websites contain at least one page in two languages (42%), 30% of crawled websites are monolingual when crawled to a depth of two links, while the rest (28%) have at least one page of content in 3 or more languages.

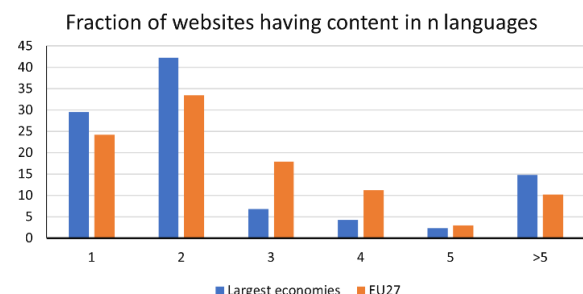


Figure 1 Fraction of websites having content in n languages (depth 2 links).

From the second list of EU 27, 401 websites were crawled and analysed (**Figure 1**). In this

¹ https://commission.europa.eu/resources-partners/etranslation_en

case, when crawled 2 links deep, one-third of websites have content in two languages, 24% of websites are monolingual and the rest (~43%) are multilingual, i.e., having at least one page of content in 3 or more languages.

As usually links to valuable content in other languages are provided on the landing page of the website, we compared these results with those from crawling with the depth of one link. In this case, less than 20% of websites from the countries with the largest economies and less than 25% of EU 27 websites yielded content in 3 or more languages. This shows that for many websites multilingual content is available only for some of the content on deeper levels.

A detailed analysis was performed on 698 public sector websites across 31 European countries (EU 27 countries plus Albania, Serbia, North Macedonia, and Switzerland) for the number of languages in which the websites are published (**Figure 2**). We found that government websites are significantly more multilingual than websites of Education (NACE (The Statistical Classification of Economic Activities in the European Community) code 85) and Healthcare (NACE codes 86, 87, 88) institutions.

3 Solutions for Website Multilinguality

Website multilingualism is enabled by an ecosystem of solutions that includes website builders, content management systems (CMS), machine translation services, systems to manage translation workflows, and plug-ins. Tools such as website builders like Elementor and WPBakery, along with eCommerce platforms like Wix and Shopify, provide a range of built-in translation tools to enable website owners to create multilingual websites without any coding experience. CMS such as WordPress, Drupal, and Joomla integrate plug-ins and extensions that allow website owners to easily translate their websites. In the backend, automated translation services such as Google Translate, Microsoft Translator and DeepL provide machine translation to automatically translate website content into multiple languages. The translation management systems (TMS) of language technology vendors like Phrase, Unbabel, Tilde, and many others also enable website translations.

3.1 User Preferences

To analyse user preferences in selecting and using automated website translation solutions, a

CATI (Computer Aided Telephone Interview) survey was conducted among 122 European companies of all sizes. The group of respondents consisted of decision makers and influencers knowledgeable of web translation topics. The countries of residence of respondents include Croatia, Estonia, France, Germany, Greece, Italy, Latvia, Lithuania, Malta, Poland, and Spain.

Our study finds that the most common reasons for businesses to translate their websites are to reach new markets (78%), improve customer service (45%), and comply with regulations (30%). The main challenges of automated website translation are accuracy (54%), cost (36%), and ease of use (32%). The most popular MT services used by

	Central government	Regional government	Local Government	Education	Healthcare
Monolingual	23%	51%	50%	48%	65%
Bilingual	50%	25%	26%	45%	26%
Trilingula	12%	12%	5%	5%	5%
4+ languages	15%	12%	19%	2%	4%

Figure 2 Multilinguality of public sector websites.

businesses are Google Translate (68%), Microsoft Translator (42%), and Amazon Translate (31%).

Smaller companies prefer easier, user-friendly, and simplified translation processes. Large companies are looking for advanced functionalities such as access based on roles, workflows that allow consistency of translations, and support for various types of content (documents, videos, blog posts, etc.).

Not surprisingly, data security and compliance are important topics for the majority of the users (71% of the users responded with "Extremely important" or "Very important" to the respective question). Security appears to be more relevant for entities with more than 100 employees than for smaller organizations. The industry sectors that care the most about the security of the websites include the Financial sector and the Distribution and Services sector.

77.9% of the users of MT solutions for translating a website would recommend the use of such solutions to other website owners. 88.4% think that it has helped improve user experience and expand their business, and 66.3% value the cost efficiency for reaching a wider audience. Among those who wouldn't recommend the usage of MT for translating websites, the reliability of translations and the quality of translated content are the major arguments against it.

3.2 Website Translation Plug-ins

Numerous multilingual website translation plug-ins offer a range of essential features to ensure effective website localization and automated translation. These features include support for translating various elements like text, images, videos, and dynamic content across posts, menus, and widgets. Integration capabilities with diverse content management systems (CMSs), eCommerce platforms, and site builders allow for seamless multilingual content creation. Automated translation workflows, often utilizing third-party machine translation services like Google Translate, Microsoft Translate, and DeepL, streamline the translation process. These plug-ins commonly incorporate automated language recognition, editors, and translation management systems (TMS) for post-editing and collaboration. Multilingual SEO support is a standard feature, enhancing visibility through URL translation, sitemaps, hreflang tags, and more. Performance-related capabilities involve cache memory and Content Delivery Networks (CDNs) to optimize website speed. Security and GDPR compliance measures are typically in place, with data encryption and access controls. Various go-to-market models include free trials, freemium versions, and subscription plans with pricing based on translation volumes and the number of supported sites. Collaboration with CMS and eCommerce platforms is a common market strategy, with plug-ins listed in partner sections on these platforms' websites.

4 Web-T Solutions

To address the need for website translation, we have developed a Web-T website translation solution. In accordance with EC requirements, the website translation solutions that are developed are free of charge, easy to use, secure, implementable on various platforms, flexible, adaptable to different CMSs, integrate free MT service eTranslation provided by the European Commission, as well as open for other machine translation providers. The key findings from the user survey and existing website translation plug-in review were included in the requirements when designing the solutions.

4.1 Overall architecture

The project solutions should suit various types of websites. The majority of websites are based on some Content Management System (CMS). Still, some websites are powered by complicated

individually built systems. On the other side of the spectrum are simpler websites that are not based on any standard or custom CMS. It should also be considered that websites can be hosted as online cloud solutions or as on-premises installations in a local hosting environment.

To cover this variety, the following types of plug-ins are being developed to reach most of the websites:

Local plug-ins developed for popular CMS platforms WordPress, Drupal, and Joomla and directly communicate with the MT service; all translations are post-edited and saved locally in the website database; machine-translation is performed in the backend, HTML page is being rendered from the local CMS database.

Universal plug-in – contains Lightweight JavaScript plug-in for any website translation; translation is performed after the page is rendered on the client's side (by the client's browser). It also includes the Translation Hub for result caching, MT provider configuration, and translation post-editing. Website translation is performed after the browser has rendered the page on the client side.

Hybrid plug-in provides a “lighter” integration in CMS platforms and encapsulates the lightweight JavaScript plug-in, which is connected to the Translation Hub. Website translation is performed after the browser has rendered the page on the client side.

Translation Hub is a distinct module designed to serve as a caching mechanism for storing and editing translations for the universal plug-in. It effectively stores content translated by the MT provider, eliminating the necessity for repetitive requests to the MT provider. Additionally, it offers a user-friendly interface for editing translations.

Each plug-in type supports two MT provider integration approaches that are implemented in the MT API Integrator:

Asynchronous eTranslation Integration – MT requests from local plug-ins and translation hubs are posted to eTranslation. The eTranslation system sends the results back to the endpoint asynchronously;

Synchronous generic MT API – generic MT API is specified. Every local plug-in and translation hub will be able to establish a connection to any MT provider that supports generic MT API implementation. This generic MT API can be created and/or hosted by the

website owners, MT providers, or any third-party translation hub host. Connection to a specific MT provider is enabled by setting the selected

Depending on the integration, the CMS local plug-in can also contain translation and language management features. For example, the WordPress/ WooCommerce plug-in has all the localisation functionality built into the plug-in, as there is no native multilingual support in the WordPress CMS. In contrast, Drupal and Joomla extensions rely on the built-in localisation features, which provide translation and language management functionality (e.g. translatable string retrieval, translation storage, editor interface, language switcher, etc.), so the main purpose of local plug-ins for Drupal and Joomla is to add automated translation functionality to the multilingual website setup.

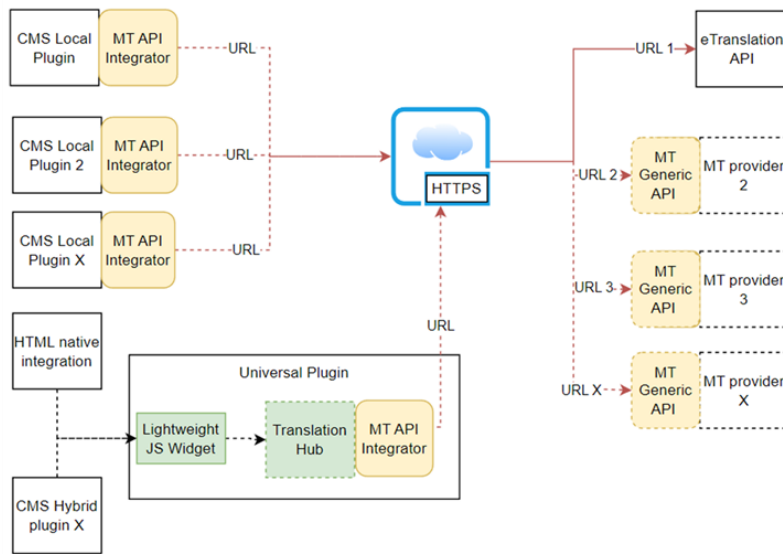


Figure 3 Conceptual architecture of Web-T multilingual plug-ins

integrator URL and an access key.

While local plug-ins will be directly downloadable from the respective CMS plug-in repositories, the universal plug-in does not have such an option, thus it will be directly downloadable from the solution website and EC code repository. Translation hubs will be hosted in a decentralised manner. Website owners can run and host the translation hub by themselves or look for any public/commercial installation available. As the translation hub is open source, any new provider can host it or create an extended solution based on it.

This architecture (**Figure 3**) provides a way to extend the WEB-T ecosystem with solutions for other CMSs without the direct involvement of all MT providers that are part of the ecosystem. The same applies to extending the ecosystem with new MT providers that will be immediately connected with all CMS integrations.

4.2 Local Plug-ins

CMS local plug-ins are installable in the respective content management system to enable the machine translation of website content. CMS local plug-ins contain an MT API integrator component. It supports asynchronous API for accessing the eTranslation service and synchronous API for communication with other MT providers.

4.3 Universal Plug-in

As the client-owned webpages can be very different in selected technology, content, and architecture, the only generic way to ensure content translation is to perform the translation after the page content and HTML are rendered. Webpages can also be interactive, so the content can also change after the initial page load has been already completed.

The rendering process is typically performed on the client-side Internet browser, thus the only reasonable technology for content translation on the client side is JavaScript code that follows the HTML content changes in end-users' browsers (**Figure 4**).

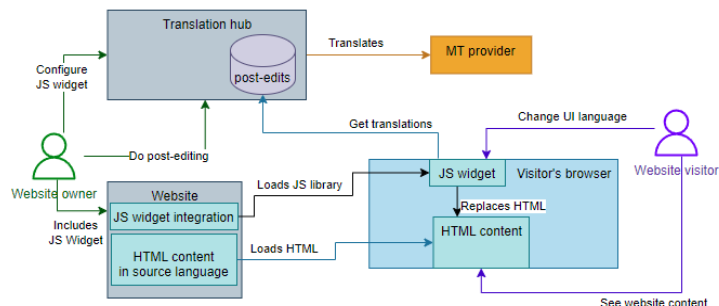


Figure 4 Conceptual architecture of Web-T universal plug-in

As JavaScript functionality is restricted and no back end is possible on the client side, an intermediate server-side tool is needed to act as a proxy between the end-user's browser and MT provider.

The necessity for the translation hub as a back-end tool arises from the need to facilitate post-editing, implement caching of MT results, secure private MT provider keys, and streamline diverse MT API workflows supported by multiple MT providers.

4.4 MT API Integrator

MT API Integrator is a specification to ensure interconnection between various CMS plug-ins and is supported by MT providers. MT API Integrator is implemented in all local plug-ins and the Translation Hub. This component consists of two parts – support for an eTranslation asynchronous approach and generic MT API for synchronous communication for the integration of any MT provider. To allow the website owner to specify which engine the MT API Integrator should use, the user interface must provide an MT engine choice in the WEB-T solution settings – eTranslation or another MT provider. With the eTranslation engine option selected, the website owner will need to provide eTranslation API credentials; when another MT provider is selected – MT provider API URL and MT provider access key.

To facilitate the integration with new CMSs, a distinct MT API integrator PHP library is created, given that PHP is the predominant language used for building CMSs.

4.5 eTranslation Integration

As eTranslation API uses digest authentication, for each call there are 2 requests – to receive authentication information and to send the actual request. Since all eTranslation API methods need authentication (including get-domains), supported language retrieval is only possible after the user has entered the valid eTranslation API credentials.

To optimise translation performance and quality using formatted text with XML or HTML tags, integrations should use document translation to send many translatable items in one request, rather than sending each string in a separate text translation request. For eTranslation integration to work, the WEB-T solution provides a REST API endpoint, which is used to receive async translation responses from eTranslation. If CMS does not support this, a local plug-in cannot be created and the hybrid approach must be used. To align asynchronous eTranslation integration workflow with other MT provider integrations (synchronous), CMS plug-ins have to wait for eTranslation responses in a synchronous way (e.g., by regularly

checking if the response has been saved in the database by the REST API endpoint handler).

5 Conclusion

Our study underscores the limited diversity of languages in the European web space and the pivotal role of automated translation tools in streamlining website localization. It highlights the need for user-friendly, accurate, and cost-effective solutions. The analysis of user requirements, the Web-T architecture, and open-source solutions offer practical guidance for extending the availability and use of automated website translation solutions. This contributes to the goal of achieving true multilinguality of European web space and advancing digital language equality in Europe.

References

- European Commission. (2008). *Communication from the Commission to the European Parliament, the Council, the European Economic and Social Committee and the Committee of the Regions - Multilingualism: an asset for Europe and a shared commitment* {SEC(2008) 2443} {SEC(2008) 2444} {SEC(2008) 2445}.
- European Commission. (2021). *The Digital Europe Programme* <https://digital-strategy.ec.europa.eu/en/activities/digital-programme>, (retrieved Oct 15, 2023).
- European Parliament. (2018). *Resolution on language equality in the digital age*. (2018/2028(INI)).
- Pastor, R., Quiros, I., Garcia, I., Cardus, I., and Nogués, M. (2017). *Language equality in the digital age-Towards a Human Language Project*. European Parliament (STOA).
- Rehm, G., Marheinecke, K., Hegele, S., Piperidis, S., Bontcheva, K., Hajic, J., Choukri, K., Vasiljevs, A., Backfried, G., Prinz, C., Gomez-Perez, J. M., Meertens, L., Lukowicz, P., van Genabith, J., Losch, A., Slusallek, P., Irgens, M., Gatellier, P., Kohler, J., Le Bars, L., Anastasiou, D., Auksoriutė, A., Bel, N., Branco, A., Budin, G., Daelemans, W., De Smedt, K., Garabik, R., Gavriilidou, M., Gromann, D., Koeva, S., Krek, S., Krstev, C., Linden, K., Magnini, B., Odijk, J., Ogrodniczuk, M., Rognvaldsson, E., Rosner, M., Pedersen, B., Skadin, I., Tadic, M., Tufis, D., Varadi, T., Vider, K., Way, A., and Yvon, F. (2020). The European language technology landscape in 2020: Language-centric and human-centric AI for cross-cultural communication in multilingual Europe. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 3322– 3332, Marseille, France. European Language Resources Association.

- Rehm, G., and Way, A. (Eds.). (2023). *European Language Equality: A Strategic Agenda for Digital Language Equality*. Springer Cham, <https://doi.org/10.1007/978-3-031-28819-7>
- Vīksna, R., Skadiņa, I., Skadiņš, R., Vasiljevs, A., and Rozis, R. (June, 2022). Assessing Multilinguality of Publicly Accessible Websites. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference* (pp. 2108-2116).
- Kelly-Holmes, H. (2019). Multilingualism and Technology: A Review of Developments in Digital Communication from Monolingualism to Idiolingualism. *Annual Review of Applied Linguistics*, 39: 24 – 39. Cambridge University Press.
- Miraz, M. H., Ali, M. and Excell, P. (2013). Multilingual Website usability analysis based on an international user survey. In *The Proceedings of the Fifth International Conference on Internet Technologies and Applications (ITA 13)*, 236-244.
- Sargent, B.B. and Lommel, A. (2019). *Global Website Assessment Index 2019*. CSA Research Report. <https://insights.csa-research.com/reportaction/48682/Marketing>.
- Wright, S. (2004). Introduction to special issue on Multilingualism on the internet. In *International Journal on Multicultural Societies* 6(1): 5-13.
- Singh, S. P., Kumar, A., Darbari, H., and Maheshwari, N. (2016). Plug-in for Instantaneous Web Page Rejuvenation and Translation. In *Smart Trends in Information Technology and Computer Communications: First International Conference, SmartCom 2016*, Jaipur, India (pp. 77-87). Springer Singapore.
- Stasimioti, M., Sosoni, V., Kermanidis, K., and Mouratidis D. (2020). Machine Translation Quality: A comparative evaluation of SMT, NMT and tailored-NMT outputs. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 441–450, Lisboa, Portugal. European Association for Machine Translation.

Appendix A Questionnaire for the interviews with End users.

Screening:

Q1: Is your company website multilingual?

If yes Q2, if no – cancel

Q2: Are you using automated translation to make your website multilingual?

If no Q3, if yes Q41

Q3: Are you using automated translation solutions for translating documents, texts or to transcribe voice communication?

If yes Q42, If no Q5

Q41: Which solution you are using?

- Register the solution.
- Not sure.

Go to Q51

Q42: Which automated translation solution are you using? (we have 4.1 as we need to distinguish between those that are used for website translations and those used for document translations)

- Register the solution.
- Go to Q51

Q5: Why are you not using automated translation of your organization's website?

Select all that apply:

- We have not tried because we need control over the accuracy of the translation and don't believe that automated translation can provide such.
- we tried, but the quality of the translation was not good enough.
- We tried but had issues such as untranslated elements of the site, layout issues, and others.
- We are not aware of what automated translation solutions may be applied for automated website translations.
- We don't have the technical skills to deal with the integration of a technical solution on our website.
- we investigated options, but the investment seems too high.

- We are faced with incompatibility with the existing IT infrastructure.

- We don't need translation of the content of the website.

- Other.

Thank you, cancel.

Q51. Where did you learn about the e-translation tools for your website?

Select all that apply:

- In social media (FB, Instagram, other).
- In specialized blog posts.
- From our Web developer.
- Found it in the online store (Shopify store).
- We researched on the internet.
- Other, please specify.

Q52. How did you choose the specific automated translation solution you are using?

Select all that apply:

- It was recommended by a colleague or friend.
- It was recommended by our web developer/IT team.
- It had the best reviews and ratings online.
- It was the most affordable option.
- It offers most of the features and functionalities that we need.
- It offers the highest quality.
- It is provided by the tools that we use to build our website (content management systems, site builder, eCommerce system).

Q6: What capabilities of an automated translation solution are important for your company with respect to the quality of translation?

Please rank from 1-5 (where 1 is not important and 5 is very important)

- To support a wide variety of languages.
- To handle industry-specific or company-specific terminology particularly well.
- To offer supreme quality of translation for specific language pairs and subject areas.

- The capability of my organization to build custom language models.
- Availability of “adaptive machine translation” models that “learn” and adapt to new words and phrases over time.
- Availability of some level of human validation of translations.
- To translate all the elements of the website – incl. widgets, product descriptions, and buttons across all web pages, custom posts, blogs.
- Others: Please specify.

Q7: Assuming that the quality of translation of an automated translation solution is good enough, what other capabilities are important when selecting a tool for automated translation of your website?

Please rank from 1-5 (where 1 is not important and 5 is very important)

- Ease of use of the solution.
- Ease of integration of the solution with the technologies used by the website.
- Availability of SEO (Search engine optimization) capabilities to improve website ranking; Examples of capabilities: translation of URLs, translation of metadata, Search Engine Indexing, (to rank higher local language searches), Search Engine Friendly (SEF) URLs i.e., dedicated URL for a multilingual SEO strategy, etc.
- Editing in context - i.e., users are able to see exactly how the translated content looks on the website.
- Post-machine translation editing capabilities allowing collaboration of different roles.
- Ability to support specific content for the language-specific versions of the site.
- Automated translation does not harm the performance of the website.
- Quality of the support from the solution provider.
- Others: Please specify.

Q8. Please, indicate to what extent the following features provided by solutions that enable the automated translation of a website are important to you:

- Please rank from 1-5 (where 1 is not important and 5 is very important)
- User interface allows to switch/cancel ad hoc the level of the service.
- Usage statistics/ dashboard.
- Data security and privacy features to prevent disclosure of confidential information.
- Ability to control access to content based on roles.
- Solution complying with GDPR, PCI, HIPAA, or other industry standards.
- Portability of the solution (ability to change the hosting provider, the provider of the CMS, etc., and to keep the vendor of the automated translation solution).
- It is possible to ask for a refund of pre-paid subscription fees.

Q9: What vendor offering options were important for selecting a translation vendor:

- Free trial
- Free version of the solution.
- Possibility to switch or cancel ad hoc the level of the service.
- Vendor policy allows to continue using the translated versions of your website if you don't renew your license.
- Hosting services, provided by the translation solution provider.
- Marketing automation capabilities built into the translation solution platform or provided by third parties.
- Affordable pricing
- Other:

Q10: Which features are missing in the current market offering of your vendor of automated translation services?

- Register1:
- Register2:
- Register3:

Optional Question:

Q11: Would you recommend the use of the automated translation services solution to other website owners? Why or why not?

- Yes, because it has helped improve user experience and expand my business.
- Yes, because it is a cost-effective solution for reaching a wider audience.
- No, the quality of automated translations is not good enough.
- No, because it is not a reliable substitute for manual translation.
- Unsure.
- Other, please specify.

Exploring the Effectiveness of LLM Domain Adaptation for Business IT Machine Translation

Johannes Eschbach-Dymanus Frank Essenberger Bianka Buschbeck Miriam Exel

SAP SE

Dietmar-Hopp-Allee 16, 69190 Walldorf, Germany

firstname.lastname@sap.com

Abstract

In this paper, we study the translation abilities of Large Language Models (LLMs) for business IT texts. We are strongly interested in domain adaptation of translation systems, which is essential for accurate and lexically appropriate translation of such texts. Among the open-source models evaluated in a zero- and few-shot setting, we find Llama-2 13B the most promising for domain-specific translation fine-tuning. We investigate the full range of adaptation techniques for LLMs: from prompting, over parameter-efficient fine-tuning to full fine-tuning, and compare to classic neural machine translation (MT) models trained internally at SAP. We provide guidance how to use training budget most effectively for different fine-tuning approaches. We observe that while LLMs can translate on-par with SAP's MT models on general domain data, it is difficult to close the gap on SAP's domain-specific data, even with extensive training and carefully curated data.

1 Introduction

With swift improvement and recent successes of Large Language Models (LLMs), it has become imperative for companies to measure their productive NLP systems against such new models. In the rapidly evolving field of NLP, incorporating the state-of-the-art models could unlock new capabilities for one's product and improve performance. On the other hand, switching to LLM-based systems

should not be done merely to appeal to public hype, but should be a thoroughly evaluated choice.

With this in mind, we set out to investigate whether LLMs can be easily utilized to outperform and ultimately supersede the current machine translation systems employed by SAP (Buschbeck et al., 2022). They are based on a traditional neural machine translation architecture trained on a multitude of data sources including the contents of the company-internal translation memories and is therefore optimized for SAP's domain of interest, which we call *Business IT* here. While previous research has shown that LLMs make good translators (Hendy et al., 2023; Zhang et al., 2023; Zhu et al., 2023), it is not yet well explored whether they can effectively adapt to domain-specific translation intricacies and outscore a smaller model that has been trained from scratch within the domain.

In particular, our interest lays in whether comparably smaller sized open-source LLMs can be fine-tuned to this end. This interest is motivated by certain drawbacks of using large proprietary models such as OpenAI's GPT-4 (OpenAI et al., 2024) out-of-the-box. Potential data privacy concerns, slower inference and higher monetary costs (provided sufficient throughput) are some of the reasons.

In addition, fine-tuning (open-source) models offers some more benefits. Fine-tuning addresses challenges such as hallucinations and overgeneration, commonly associated with the LLMs' innate generative nature. By channelling the LLMs' focus towards translation through downstreaming, it becomes possible to regulate and control these unwanted generative tendencies, resulting in a more precise and tailored output for the intended domain-specific application.

While there is a general argument to be made that fine-tuning an LLM on parallel data would im-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

prove its translation quality, a much more crucial one can be made regarding domain-specific translation. In domain-specific translation, vocabulary and translation patterns might differ substantially from its general purpose counterpart. For instance, UI strings contained in our domain-specific data such as ‘Item masters in catalog’, ‘Number in the Total field’ or ‘PENDING The item has not yet been sent.’ do not only differ from general domain texts in terms of vocabulary, but can also deviate substantially in syntax.

In our experiments, we first estimate the translation capacities of various models through prompting to identify the most promising one to fine-tune. Concretely, Llama-2 13B appears to offer the best model size to performance trade-off. We then fine-tune the model with both full parameter and parameter-efficient tuning and conduct various ablation studies.

Overall, we arrive at the following conclusions:

1. It requires large amounts of domain-specific data to have a fine-tuned Llama-2 13B approach the performance of a smaller but dedicated translation system.
2. When fine-tuning Llama-2 13B with domain-specific translation data, we do not observe noteworthy catastrophic forgetting of general domain translation capacity.
3. For domain-specific translation, low-rank adaptation (Hu et al., 2021) cannot compete with full fine-tuning as it fails to internalize domain-specific phenomena in the limited parameters available. In order to increase model fit, the adapter rank needs to be increased to magnitudes where the tuning becomes no longer parameter-efficient. Additionally, we observe that quantization significantly hinders model fit and consequently in-domain performance as well.
4. If nonetheless parameter efficient fine-tuning is conducted, one should favour an increase in training data over an increase of training epochs. For full fine-tuning, however, training for multiple epochs provides a notable benefit. In fact, if the number of training iterations are to be kept a constant, one should favour an increase of epochs over more training data; naturally, under diminishing returns.

2 Related Work

In their paper, Brown et al. (2020) presented the performance of GPT-3 and evaluated the translation ability of their LLM. They found that general-purpose LLMs benefit from having examples in the prompt (few-shot prompting) to guide the model towards a specific task. Undoubtedly, adding relevant examples via few-shot prompting or retrieval augmented generation can improve translation performance. However, both Alves et al. (2023) and Li et al. (2023) observe that translation fine-tuning outperforms few-shot prompting when provided with only few thousands of training samples. Xu et al. (2023) achieve state-of-the-art translation performance with help of a two-stage training mechanism, where in the first stage, the model is further pre-trained on billions of tokens of monolingual data of various languages to shift the model to a more multilingually balanced state; away from its dominantly English pre-trained state. Only then, the model is fine-tuned with limited parallel data. While the resulting performance is astonishing, the first training stage is computationally expensive. Even with adequate GPU resources available, the proposed setup does not necessarily work as effective in low-resource domain translation, where the parallel data does not align well with the monolingual data used in the first stage. Üstün et al. (2024) recently presented the Aya model which uses a more balanced distribution of multilingual data in the pre-training. Although the approach seems promising, our first preliminary investigations do not show substantially higher translation performance of the Aya model compared to previous LLMs.

3 Datasets & Evaluation

In this paper, we mainly focus on the high-resource language pair *English* \rightarrow *French*. The more resources a language pair has, the more LLMs should be able to leverage from their pre-training, making it easier and quicker to downstream them for the translation task. In addition, we also investigate the performance on the low-resource language pair *English* \rightarrow *Slovak* and on *English* \rightarrow *Japanese*, which is known for its complexity, in section 5.7.

For few-shot example retrieval, fine-tuning and testing, we use well-curated parallel SAP-internal data. It is composed of large amounts of software user interface (UI) strings, user assistance (UA) texts, but also training materials, corporate content and marketing texts. The models are tested not

only on a test set of 2000 segments consisting of domain-specific UI strings and UA texts, but also on general-domain data, i.e. the FLORES (Goyal et al., 2022) test set. Even though they stem from the same domain (SAP), training and in-domain test data are not merely divided randomly, but instead feature both temporal and distributional shifts. This in turn allows a more realistic performance evaluation (Søgaard et al., 2021).

We want to clarify that SAP’s MT systems are trained using large corpora comprised of millions of parallel sentences, and this training is performed over many epochs. Since fine-tuning LLMs on a similar scale would entail considerable computational costs, we conduct our fine-tuning experiments with fewer but gradually incremented quantities of data to map respective improvements in translation quality. Furthermore, while SAP’s MT system evaluated in this study has been trained to excel in translating texts from the SAP domain, it has not been fine-tuned to UI and UA texts specifically, and obviously the test data is unseen.

We evaluate the performance with both BLEU (Papineni et al., 2002) and COMET¹(Rei et al., 2022). While COMET is more robust and correlates better with human annotators, the n-gram based BLEU score nonetheless has its use when evaluating domain-specific translations. Specifically, it captures lexical agreement with references which indicates the correct use of terminology and writing style for the domain.

4 Prompting

To establish a baseline for the translation fine-tuning it is natural to start with simple prompt experiments. These experiments are relatively straightforward to conduct since large commercial models like GPT-4 are offered as services. Furthermore, hosting open-source models for inference requires fewer resources compared to tuning them. The motivation driving this inquiry is twofold: first, to establish a baseline for the performance of open-source models; and second, to evaluate the inherent capabilities of GPT-4, a proprietary cutting-edge language model in its “out-of-the-box” state.

Given the nature of the task, we restrict ourselves to a selection of models that are intended for multilingual usage:

1. **GPT-4:** GPT-4 serves as a benchmark for the

¹<https://huggingface.co/Unbabel/wmt22-comet-da>

state-of-the-art in natural language processing, boasting superior language understanding and generation capabilities (OpenAI et al., 2024). Hendy et al. (2023) have also demonstrated that it shows remarkable performance in the translation task.

2. **Llama-2 (7B, 13B, 70B):** The Llama-2 family of models has been shown to achieve great performance in various tasks across various languages and has a commercially usable licence.
3. **BLOOM 7B:** The BLOOM model family has been released in 2022 and was trained on well documented high-quality data (Laurençon et al., 2022) encompassing 46 natural languages.
4. **Falcon (7B, 40B):** As with BLOOM, the Falcon family of models, released in 2023, is of special interest due to its balanced, curated and, most importantly, well documented multilingual training data (Penedo et al., 2023).

English: Receive supplier invoice
French:

Figure 1: Translation Prompt

The simple prompt shown in figure 1 is used for prompting and fine-tuning experiments throughout the paper. While it is known that optimized and more verbose prompts can improve results, we refrain from *prompt engineering* for two reasons. Firstly, while engineering prompts is fairly cheap when optimizing a zero-shot setting, it would require repeated trainings for each and every prompt to measure its performance in a fine-tuning setting; an endeavour that is too costly. In our experiments, we expect the model to adapt to any prompt during fine-tuning. Susceptibility to prompt design would prove a major obstacle for fine-tuning LLMs. Secondly, a short and concise prompt is preferable as it leaves more context length available for the actual translation pairs.

Figure 2 displays the BLEU and COMET scores of the models in the zero-shot setup. For open-source models, we establish a comparison under equal resource conditions, i.e. given four NVIDIA A10G². Models that are too large for full precision inference are tested with 8-bit quantization instead.

As expected, SAP’s MT model performs best by a large margin on the domain-specific data. The superior BLEU score, in particular, indicates the cor-

²a single *g5.12xlarge* instance on AWS

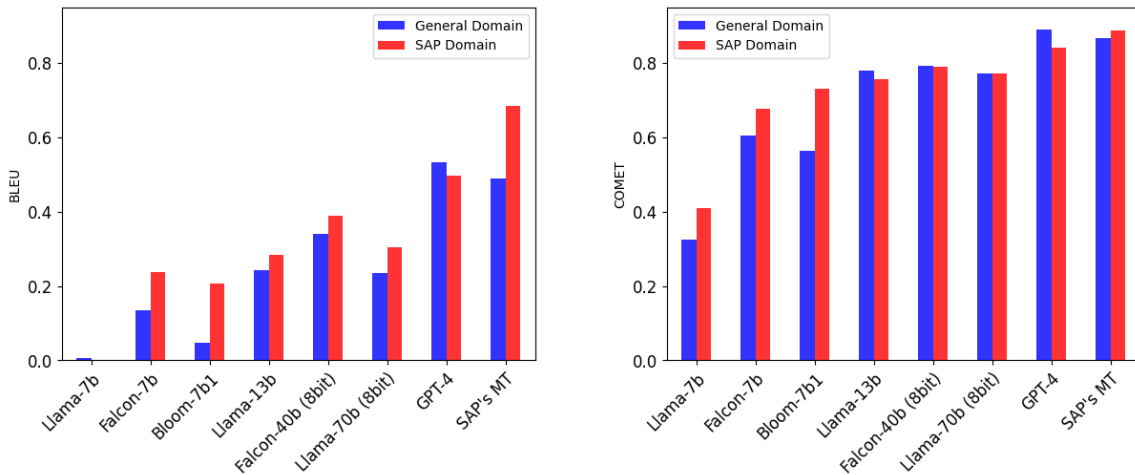


Figure 2: Zero-shot prompting performance

rect usage of in-domain vocabulary, rather than just semantically similar phrases. For general-domain data (FLORES), however, GPT-4 is capable of outperforming SAP’s domain-specific MT system out of the box. Keep in mind that this MT model is not optimized for general-domain translation, and it is unknown whether the publicly available dataset was included in GPT-4’s training data.

Provided with the simple prompt, all open-source models but Llama-2 7B are able to translate texts consistently³, albeit not necessarily correctly. Naturally, it is difficult to determine whether the weaker model performance is the result of shortcomings in its translation quality, or merely a misinterpretation or mishandling of the prompt.

It might seem unintuitive, but some open-source models perform better on domain-specific data than on general-domain data. This can be explained by the large percentage of UI segments contained in the SAP test set. While such strings contain lexical intricacies, they are generally short and syntactically simple, which makes them easier to translate for the smaller models.

As a natural next step, we investigate whether few-shot prompting could narrow the gap for the SAP domain-specific translations. We focus our experiments on GPT-4, serving as an upper LLM benchmark, and Llama-2 13B, which offered the best trade-off between model size and translation performance in the zero-shot experiments. In order to construct our few-shot prompts, we first encode the English source segments of the domain-specific parallel training data (section 3) with the sentence-BERT model *all-MiniLM-L6-v2* by Reimers and

³generating output in French and not English

Gurevych (2019). Then, for each English segment to translate, we retrieve the five translation pairs that have the highest cosine similarity to the English source. These pairs are then arranged as completed prompts and placed before the final segment set for translation, with each pair separated by an empty line.

For the Llama-2 13B model we had to conduct postprocessing of the output due to overgeneration issues. In particular, it continued generating English-French sentence pairs beyond the completed prompt. To deal with this, we simply truncated the generation after the first line break. The results of the few-shot experiments are displayed in table 1. While both models improve in performance of in-domain translation, a gap to the domain-specific MT model still remains. However, it is worth noting that the open-source Llama-2 model benefits more substantially from the examples, promising potential that may be even better leveraged through fine-tuning. We also find that providing domain-specific examples for general-domain translation is detrimental to the models’ performance.

While few-shot prompting can provide a basic understanding of the domain in question, it may not fully capture all the domain-specific nuances. This is particularly challenging when the domain is highly specialized, as selecting the appropriate domain-specific vocabulary and creating accurate examples can be difficult. Additionally, including multiple examples in the prompt increases the token count, which can lead to higher computational costs and longer inference times. Fine-tuning, on the other hand, would allow the model to learn

Model	SAP domain	General domain
Llama-2 0-Shot	0.757 / 0.283	0.780 / 0.243
Llama-2 5-Shot	0.843 / 0.503	0.767 / 0.310
GPT-4 0-Shot	0.842 / 0.498	0.891 / 0.532
GPT-4 5-Shot	0.866 / 0.560	0.883 / 0.518
SAP’s MT	0.888 / 0.686	0.867 / 0.490

Table 1: Few-shot results with COMET (first value) and BLEU (second value) for Llama-2 13B and GPT-4. For comparison SAP’s MT is added.

from a much larger pool of examples, potentially leading to better adaptation to the specific-domain requirements.

5 Fine-Tuning

While the performance of open-source LLMs is promising, it falls short compared to both a domain-specific neural machine translation model as frequently used in production and more advanced models like GPT-4. A promising pathway to improve the task-specific performance of open-source LLMs is to further fine-tune them. Given the substantial gap in performance on the in-domain data to SAP’s MT model and even GPT-4, we have a strong incentive to investigate how well LLM translation can adapt to a specific domain through fine-tuning. To do so, we experiment with three different fine-tuning setups:

1. **LoRA:** With low rank adaptation (Hu et al., 2021), the pre-trained model weights are frozen while trainable rank decomposition matrices are injected on top of the frozen weight matrices. As the decomposition matrices are the only ones fine-tuned and contain magnitudes less parameters, model downstreaming becomes faster and more GPU efficient.
2. **QLoRA:** The fine-tuning approach proposed by Dettmers et al. (2023) quantizes the pre-trained model during training and only keeps the trainable LoRA adapter weights in standard precision. This method reduces the memory requirements of fine-tuning which in turn, depending on the available GPUs and model size, can allow data-parallel training rather than model-distributed one, cutting training time short by a multitude.
3. **Full fine-tuning:** While Hu et al. (2021) and Dettmers et al. (2023) show that the proposed parameter efficient fine-tuning approaches perform on-par with full fine-tuning, other researchers applying them could not always con-

firm such observations (Sun et al., 2023; Chen et al., 2022). Consequently, we also conduct full fine-tuning to establish an upper bound.

QLoRA is of special interest, as the greatly reduced GPU footprint allows cost-efficient training. In addition, quantization could also help reduce the cost during inference and recent development in dynamic adaptation (Babakniya et al., 2023) make QLoRA even more tempting. A main interest of our experiments is therefore an evaluation of QLoRA against full fine-tuning for domain-specific translation. Then, ablation studies are conducted that investigate shortcomings of QLoRA opposed to LoRA without quantization.

5.1 Fine-Tuning Setup

We use the 13 billion parameter version of Llama-2 for all fine-tuning experiments. For one, the model is capable of translation in a zero-shot prompt setup, which certifies that there is sufficient pre-trained knowledge to leverage through fine-tuning and allows a comparison to a sensible baseline. Secondly, the model is comparably lightweight, which allows full fine-tuning on as few as four NVIDIA A10G⁴. We use standard libraries to perform the fine-tuning, namely Huggingface’s trainer interface⁵ and bitsandbytes⁶ for quantization. We use the training data presented in section 3 and vary the amount of training segments in the experiments and train for 3 epochs.

Measuring performance not only in the domain-specific but also on general domain data allows us to investigate the effect the domain-specific translation tuning has on the model’s translation performance in general. On the one hand, we expect an increase in general domain translation performance, as the model is downstreamed to translate only. On the other hand, increasing the model fit to specific data could also induce catastrophic forgetting and consequently cause the general domain performance to deteriorate.

For (Q)LoRA training, we set the rank to $r = 8$ and the scaling factor to $\alpha = 16$ unless otherwise specified. We use 8-bit quantization for QLoRA, as its 4-bit counterpart did not yield satisfying results in preliminary experiments. For both (Q)LoRA and full fine-tuning, we observed good convergence

⁴Using paged optimizers as discussed in Dettmers et al. (2023)

⁵https://huggingface.co/docs/transformers/main_classes/trainer

⁶<https://github.com/TimDettmers/bitsandbytes>

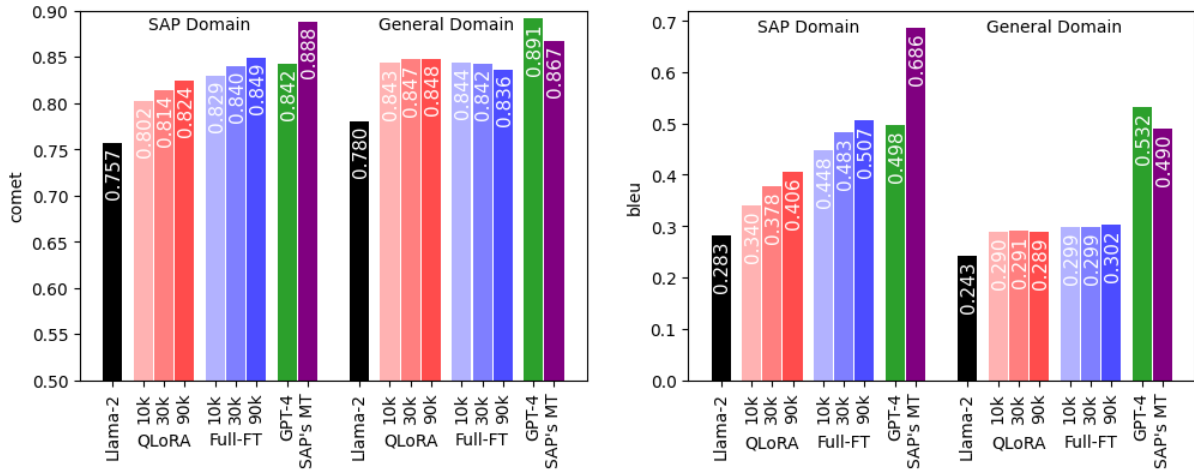


Figure 3: Model performance measured in COMET (left) and BLEU (right). The Llama-2 models have been fine-tuned for 3 epochs with 10 to 90 thousand parallel segments.

behaviour with a learning rate of $2e-5$. Other than heuristically searching for a functional learning rate, we did not search any further hyperparameters. After all, one key advantage of (Q)LoRA over full fine-tuning is that it is not as sensitive to hyperparameters; an advantage we do not want to offset by expending valuable resources to optimize full-fine-tuning.

5.2 Results

The results in figure 3 demonstrate how effective LLMs can learn from very limited training data. Even a small training set of only 10k sentence pairs drastically improves performance over the zero-shot baseline. This initial boost of performance, compared to zero-shot, is most likely due to the model quickly adjusting to the prompt and translation task in general. By increasing the training data we can further improve the model’s performance, albeit with diminishing returns. At 90k training samples, the fine-tuned model surpasses GPT-4 performance on the domain-specific test sets. This demonstrates that fine-tuning of a smaller open-source model can close the performance gap to large proprietary models out of the box. With limited training data, however, Llama-2 cannot be easily downstreamed to beat the parameter efficient SAP translation system. Further investigations into the amount of training data required to match SAP’s MT performance are conducted in section 5.6.

Despite training only on domain-specific SAP data, the model also shows improvements in general domain translation performance. While this is unsurprising, given that the model is downstreamed

for translation, it is nonetheless remarkable that there is no apparent catastrophic forgetting occurring when fine-tuning with the above quantities of training data. When full fine-tuning, we begin to see a slight degradation of general domain performance from 30k samples upwards. However, the performance is still substantially better than the model’s zero-shot one. This general robustness also stands in contrast with few-shot prompting, where the addition of domain specific examples deteriorates general domain performance. In a way, one could argue that the few-shot examples much more aggressively urge the model to translate in the domain-specific style while fine-tuning only provides the models with the additional knowledge to translate appropriately, if necessary.

In general, we observe that full fine-tuning is superior to QLoRA tuning. Most notably, however, is that the full fine-tuned model displays a much larger improvement in BLEU scores on domain-specific data than its QLoRA counter-part. Since BLEU is a token-based metric, we conclude that full fine-tuning allows the model to internalize the lexical intricacies of the domain. This is crucial for translation use cases in specific domains, also at SAP, as the translations must be consistent with established terminology. With less trainable parameters available, QLoRA is less capable of internalising these lexical differences.

While full fine-tuning is undoubtedly the superior choice, it comes with increased computational costs. With the GPU setup discussed above, the QLoRA training could be conducted in a data distributed manner, while the full fine-tuning required

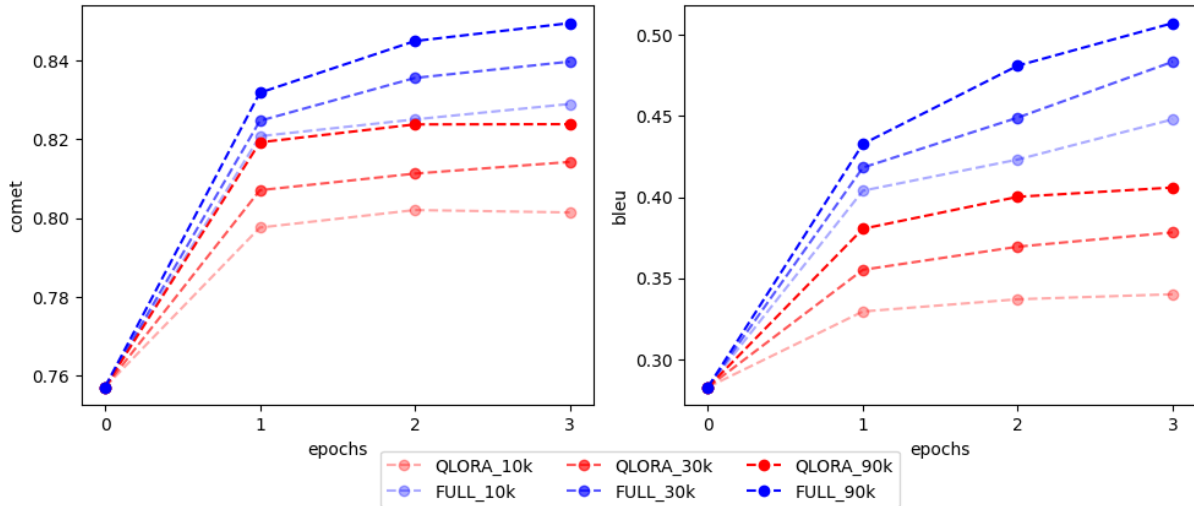


Figure 4: Effect of increasing training data and number of epochs on BLEU and COMET scores for the SAP domain test set.

all four GPUs for a single copy of the model.

5.3 Budget Efficient Training

When confronted with a fixed training budget, the pivotal decision arises between allocating resources to acquiring more training data or investing in multiple epochs. Iterating over the same data across multiple epochs expedites model fitting but poses the potential threat of overfitting, as the model may become too closely tailored to the training set. Conversely, augmenting the volume of training data holds the promise of bolstering the model’s generalization capacity, yet it introduces the risk of underfitting.

Since full fine-tuning and QLoRA differ substantially in terms of trainable parameters and therefore also in expected time to reach training convergence, we investigate the effect of data and epoch increase on model performance. Figure 4 shows that QLoRA fine-tuning does not benefit all too much from an increase of epochs. After around two epochs, the test set performance is already saturated. More importantly, however, increasing the amount of training data by a factor of three provides a substantially larger boost in performance than tripling the training epochs. Therefore, we conclude that when fine-tuning with QLoRA it is sufficient to have the model observe each example only once. A reason to this is likely that the model cannot fit individual examples arbitrarily well as both the base model precision and the underparameterised adapters regulate model fit. The only way to increase the performance of QLoRA tuned models is therefore to increase the amount of training data

to allow the model to capture the underlying data distribution more wholly.

With full fine-tuning, on the other hand, we can observe a clear benefit when tuning the model over multiple epochs. Here, tripling the number of epochs results in equal or better performance than increasing the training data by a factor of three. With more trainable parameters, the model can improve its fit on an individual example with each visit. Naturally, increasing the epochs further and further will result in diminishing returns in terms of test set performance or might even lead to overfitting. Nonetheless, if confronted with limited computational budget, one should consider reducing the training data in favour of more than one epoch of training.

We would like to emphasize how differently the LLMs learn compared to the traditional and much smaller encoder-decoder translation models. These models are trained with substantially more data over dozens of epochs, since training is much cheaper, faster and requires comparably few parameters. In contrast, we see that the LLMs are quickly and easily adjusted to a downstream task in a few epochs and with a few thousand examples. This compensates for the higher training cost per sample due to the large model size.

5.4 Domain-Specific Translation - Appetite for Parallel Data

For both full fine-tuning and QLoRA, the test set performance continuously increases logarithmically with respect to the amount of training data. This observation stands in firm contrast to Xu et al. (2023),

who notice a lack of improvement beyond 10k translation examples when fine-tuning a Llama-2 7B model. Consequently, they argue that LLMs are not hungry for parallel data and suffer from catastrophic forgetting (French, 1999; Kirkpatrick et al., 2016) when confronted with too many examples. These conflicting observations could, for one, be explained by the larger size of the Llama-2 model employed in our experiments. After all, robustness to catastrophic forgetting scales with model size (Dyer et al., 2022).

We argue, however, that the different observations could stem from the type of training and test data, rather than the models. Xu et al. (2023) tune their model on general domain translation data, for which the LLM already contains knowledge that can be leveraged. Consequently, the fine-tuning just needs to nudge the model in the right direction to utilize this intrinsic knowledge.

In our case, however, we fine-tune the model on domain-specific data, which poses two challenges for the model. For one, the model needs to ‘learn’ the domain to retrieve relevant pre-trained knowledge. With an increase of training data, the model can come to a better understanding of what the domain really entails. Much more, however, the model also needs to internalize very rare or even new information it encounters during training. The update in parameters required for this in turn is much larger than the small nudge required for general domain translation.

5.5 Tackling Shortcomings of (Q)LoRA

The experimental results show a substantial gap in performance between QLoRA tuning and full fine-tuning. To attempt to close this gap, we experimented with various configurations for QLoRA training.

We note three observations based on the results in table 2. First, while applying adapters to all attention and feed forward matrices provides a substantial boost, it is still not comparable to full fine-tuning. Second, making the LoRA adapter bias terms trainable does not yield any benefit.

Finally, since the QLoRA adapters are of low rank, our suspicion was that the number of parameters is simply insufficient to learn the domain-specific intricacies of the translation data. Intuitively, with the rank approaching the full rank of the matrix and applying it to all matrices, we should also observe the performance converge towards the

Adapt. Attention	Adapt. FFN	QLoRA Bias	QLoRA Rank	BLEU
✓			8	0.406
✓			32	0.408
✓			64	0.407
	✓		8	0.381
✓	✓		8	0.458
✓	✓		64	0.456
✓	✓	✓	8	0.456

Table 2: BLEU scores on the in-domain test data for different QLoRA configurations fine-tuned with 90k parallel segments. *Adapt. Attention* signifies the low rank adaptation of the model’s *query, key, value* and *out* projection matrices within the attention submodule. *Adapt. FFN* signifies the low rank adaptation of the *up, gate* and *down* matrices within the model’s MLP submodule. *QLoRA Bias* indicates whether the low-rank adapter contains tunable bias terms. *QLoRA Rank* specifies the rank of the low-rank approximation matrices.

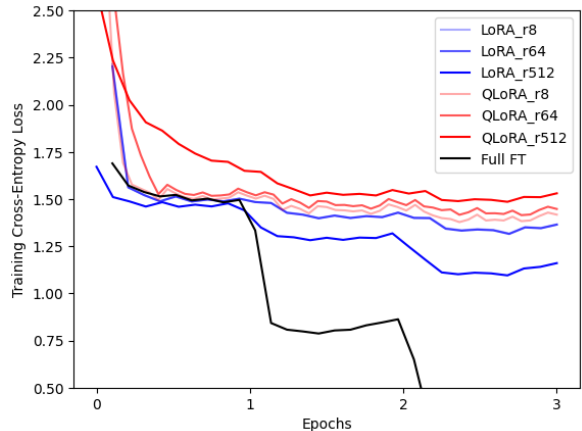


Figure 5: Effect of quantization and adapter rank on model fit.

one of the full fine-tuning. However, increasing the rank does not result in any improvement. Even worse, training runs with even higher ranks resulted in continuously degrading performance. A small grid search over learning rates and LoRA α terms could not alleviate this issue.

As figure 5 shows, QLoRA training runs are underfitting, converging quickly to a loss of about 1.5 and only very slowly beyond, regardless of adapter rank. Full fine-tuning, on the other hand, is able to fit the data much better, which becomes apparent in the dips the loss curve takes with each epoch. After all, the more often an example has been visited, the smaller the loss on it in future iterations. While the full fine-tuning model’s loss curve suggests overfitting, this is not the case yet after three epochs and validation scores are substantially better than the LoRA and QLoRA runs.

Since increasing the rank of QLoRA adapters does not result in similar fitting behaviour, we cannot hold the number of parameters alone responsible for the bad model fit. Consequently, we

investigated quantization as possible culprit. For non-quantized LoRA, figure 5 shows indeed a positive correlation between adapter rank and model fit, confirming this suspicion. With sufficiently large rank, we can once again observe the desired loss-dip at epoch boundaries, indicating that the model can fit individual examples rather than just the translation task in general. This is verified by the COMET/BLEU scores starting to converge towards the full fine-tuning scores.

A hypothesis to why quantization acts as a bottleneck might be that it causes the model to lose fine-grained information. While this loss might not be apparent when prompting the model, it could become noticeable during fine-tuning. After all, fine-tuning the model allows us to better leverage the relevant pre-trained knowledge that could not be accessed as easily through prompting. If this knowledge in turn is encoded in higher precision variations in the parameters, quantization would inevitably result in its loss.

While keeping the base model unquantized improves adapter-based fine-tuning, the BLEU scores still lack behind full fine-tuning. In order to approach full fine-tuning performance, one has to increase the rank beyond 512, which negates the advantages that *low* rank adaptation would offer in the first place. The computational load for such high ranks is comparable to a full fine-tuning, since the hidden layer size for Llama-2 13B is 5120.

5.6 Pushing The Limits

Up to now we have investigated parameter efficient techniques and compared them to a full fine-tuning. The results indicate that only full fine-tuning of LLMs can possibly lead to results comparable to GPT-4 and the encoder-decoder MT system used at SAP. Therefore, we conducted a full fine-tuning with larger datasets, 200k and 400k, to push the limits. Due to a lack of improvement in general domain translation (see figure 3), we upsampled non-UI texts to diminish the dominance of simple and short UI strings. These non-UI texts are closer to general domain translation, featuring syntactically complete sentences rather than just phrases. With this change in the data mixture, we hope to see further fine-tuning improvements on both test sets.

Figure 6 shows that the domain-specific performance of the Llama-2 model approaches the one of SAP’s MT with increasing amount of training

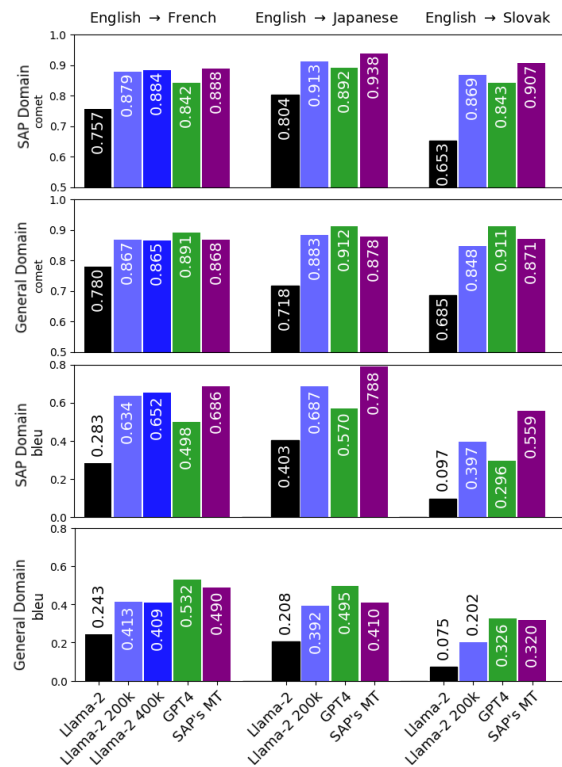


Figure 6: Llama-2 performance for larger training data sets and additional language pairs.

data. Further increasing the amount of training data will likely allow the fine-tuned LLM to surpass the MT system on the domain-specific test set. However, the same cannot be stated for general domain performance. While increasing the proportion of non-UI training data helped exceed the general domain performance observed in figure 3, we observe that further doubling the total training data does not lead to further improvement.

Possibly, further data balancing and increases in training data could allow us to fine-tune Llama-2 to match SAP’s MT on both test sets. Nonetheless, it becomes apparent that downstreaming an open-source LLM to outperform a smaller dedicated translation model is **no trivial task**.

5.7 Additional Language Pairs

To complete the picture, we also conducted experiments with two additional language pairs: English → Japanese, known for its complexity, and English → Slovak, a low-resource language pair. The results are also presented in figure 6. They show the same trends as for English → French. For the general domain, the performance saturated quickly in the same regime as the SAP MT system. For Slovak and Japanese, GPT-4 performs best on the general

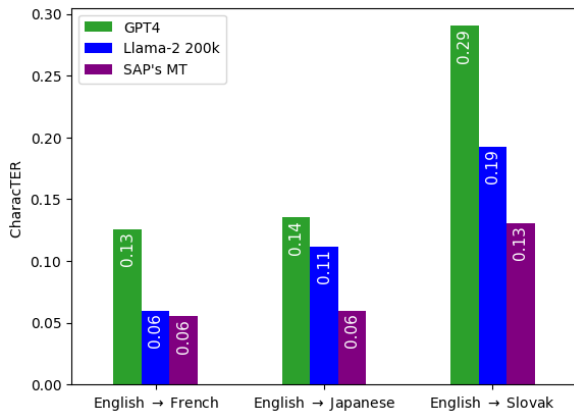


Figure 7: Human Evaluation.

domain, but as discussed before, GPT-4 has potentially seen the general domain data in training. For the SAP domain, the Llama-2 model surpasses GPT-4 and approaches the performance of the SAP’s MT system with sufficient training data.

5.8 Human Evaluation

To validate the automatic scores, we conducted a human evaluation on 300 sentences randomly selected from the SAP test set for all three language pairs. The translations generated by GPT-4, Llama-2 200K, and SAP’s MT were post-edited by two professional translators familiar with the SAP domain. We used CharacTER (Wang et al., 2016) to calculate the edit distance between the MT output and the post-edited version, and averaged the results from both translators. A lower edit distance suggests a higher quality of translation. The results, as shown in figure 7, largely corroborate the automatic metrics reported in figure 6.

6 Conclusion

We have shown that Llama-2 13B shows great potential for domain-specific translation fine-tuning and can substantially improve over its zero-shot performance. However, doing so is no trivial task. We find that few-shot prompting is not sufficient to close the performance gap to productive systems in the SAP domain. Parameter efficient fine-tuning with low rank adaptation fails to internalize domain-specific phenomena and therefore cannot compete with a full fine-tuning. Full fine-tuning, however, requires substantially more GPU compute power, which in turn is reflected in increased monetary costs.

While we were able to approach the performance of the comparably small encoder-decoder MT sys-

tem trained and employed at SAP by continuously increasing training data, we were unable to surpass it. Considering the much higher monetary inference costs and lower inference speed of the LLM compared to the MT model, the benefit of switching systems is not immediately obvious, especially when separate models would be hosted for various languages.

It is without doubt, however, that with rapidly improving released open-source models the performance for domain-specific LLM translation fine-tuning is bound to increase as well. Therefore, a continuous investigation into the translation capabilities of future open-source models is imperative.

7 Limitations and Future Work

Since LoRA is not sufficient to close the gap to full fine-tuning, we believe that multilingual fine-tuning could be a way to achieve better performance. This approach would also be more cost-efficient than fine-tuning LLMs individually for each language pair, considering their large parameter size. The emergence of multilingual models like Üstün et al. (2024) or Alves et al. (2024) makes this route even more promising. Particularly, the Tower model, which is based on Llama-2, seems to be a promising candidate. We plan to conduct multilingual fine-tuning experiments with this model in the future.

Finally, it should be pointed out that in our study, we only fine-tuned and evaluated the models at the sentence level. However, with their large context windows, LLMs are not limited to sentence-level translations and could translate whole documents. This could be especially beneficial in the SAP domain, where the consistent translation of whole technical documents is important. Progress might be more promising with fine-tuning and inference at the document level.

References

- Alves, Duarte, Nuno Guerreiro, João Alves, José Pombal, Ricardo Rei, José de Souza, Pierre Colombo, and Andre Martins. 2023. Steering large language models for machine translation with finetuning and in-context learning. In Bouamor, Houda, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11127–11148, Singapore, December. Association for Computational Linguistics.
- Alves, Duarte M., José Pombal, Nuno M. Guerreiro, Pedro H. Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal,

- Pierre Colombo, José G. C. de Souza, and André F. T. Martins. 2024. Tower: An open multilingual large language model for translation-related tasks.
- Babakniya, Sara, Ahmed Roushdy Elkordy, Yahya H. Ezzeldin, Qingfeng Liu, Kee-Bong Song, Mostafa El-Khamy, and Salman Avestimehr. 2023. Slora: Federated parameter efficient fine-tuning of language models.
- Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Larochelle, H., M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Buschbeck, Bianka, Jennifer Mell, Miriam Exel, and Matthias Huck. 2022. “Hi, how can I help you?” Improving machine translation of conversational content in a business context. In Moniz, Helena, Lieve Macken, Andrew Rufener, Loïc Barrault, Marta R. Costa-jussà, Christophe Declercq, Maarit Koponen, Ellie Kemp, Spyridon Pilos, Mikel L. Forcada, Carolina Scarton, Joachim Van den Bogaert, Joke Daems, Arda Tezcan, Bram Vanroy, and Margot Fonteyne, editors, *Proceedings of the 23rd Annual Conference of the European Association for Machine Translation*, pages 191–200, Ghent, Belgium, June. European Association for Machine Translation.
- Chen, Guanzheng, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. 2022. Revisiting parameter-efficient tuning: Are we really there yet? In Goldberg, Yoav, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2612–2626, Abu Dhabi, United Arab Emirates, December. Association for Computational Linguistics.
- Dettmers, Tim, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. QLoRA: Efficient finetuning of quantized LLMs.
- Dyer, Ethan, Aitor Lewkowycz, and Vinay Ramasesh. 2022. Effect of scale on catastrophic forgetting in neural networks. In *ICLR*.
- French, Robert M. 1999. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3(4):128–135.
- Goyal, Naman, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.
- Hendy, Amr, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are GPT models at machine translation? A comprehensive evaluation.
- Hu, Edward J., Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *CoRR*, abs/2106.09685.
- Kirkpatrick, James, Razvan Pascanu, Neil C. Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. 2016. Overcoming catastrophic forgetting in neural networks. *CoRR*, abs/1612.00796.
- Laurençon, Hugo, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, Jörg Froberg, Mario Šaško, Quentin Lhoest, Angelina McMillan-Major, Gérard Dupont, Stella Biderman, Anna Rogers, Loubna Ben allal, Francesco De Toni, Giada Pistilli, Olivier Nguyen, Somaieh Nikpoor, Maraim Masoud, Pierre Colombo, Javier de la Rosa, Paulo Villegas, Tristan Thrush, Shayne Longpre, Sebastian Nagel, Leon Weber, Manuel Romero Muñoz, Jian Zhu, Daniel Van Strien, Zaid Alyafeai, Khalid Almubarak, Vu Minh Chien, Itziar Gonzalez-Dios, Aitor Soroa, Kyle Lo, Manan Dey, Pedro Ortiz Suarez, Aaron Gokaslan, Shamik Bose, David Ifeoluwa Adelani, Long Phan, Hieu Tran, Ian Yu, Suhas Pai, Jenny Chim, Violette Lepercq, Suzana Ilic, Margaret Mitchell, Sasha Luccioni, and Yacine Jernite. 2022. The bigscience ROOTS corpus: A 1.6TB composite multilingual dataset. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Li, Jiahuan, Hao Zhou, Shujian Huang, Shanbo Cheng, and Jiajun Chen. 2023. Eliciting the translation ability of large language models via multilingual finetuning with translation instructions.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang,

- Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giam Battista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorný, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Shephard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. GPT-4 technical report.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Weijing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Isabelle, Pierre, Eugene Charniak, and Dekang Lin, editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Penedo, Guilherme, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: Outperforming curated corpora with web data, and web data only.
- Rei, Ricardo, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névéal, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Reimers, Nils and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11.
- Søgaard, Anders, Sebastian Ebert, Jasmijn Bastings, and Katja Filippova. 2021. We need to talk about random splits. In Merlo, Paola, Jorg Tiedemann, and Reut Tsarfaty, editors, *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1823–1832, Online, April. Association for Computational Linguistics.

- Sun, Xianghui, Yunjie Ji, Baochang Ma, and Xianggang Li. 2023. A comparative study between full-parameter and lora-based fine-tuning on chinese instruction data for instruction following large language model.
- Wang, Weiyue, Jan-Thorsten Peter, Hendrik Rosendahl, and Hermann Ney. 2016. CharacTer: Translation edit rate on character level. In Bojar, Ondřej, Christian Buck, Rajen Chatterjee, Christian Federmann, Liane Guillou, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Aurélie Névéal, Mariana Neves, Pavel Pecina, Martin Popel, Philipp Koehn, Christof Monz, Matteo Negri, Matt Post, Lucia Specia, Karin Verspoor, Jörg Tiedemann, and Marco Turchi, editors, *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 505–510, Berlin, Germany, August. Association for Computational Linguistics.
- Xu, Haoran, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models.
- Zhang, Xuan, Navid Rajabi, Kevin Duh, and Philipp Koehn. 2023. Machine translation with large language models: Prompting, few-shot learning, and fine-tuning with QLoRA. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 468–481, Singapore, December. Association for Computational Linguistics.
- Zhu, Wenhao, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. Multilingual machine translation with large language models: Empirical results and analysis.
- Üstün, Ahmet, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. 2024. Aya model: An instruction finetuned open-access multilingual language model.

Creating and Evaluating a Multilingual Corpus of UN General Assembly Debates

Hannah Béchara
Hertie School
Berlin, Germany
bechara
@hertie-school.org

Krishnamoorthy Manohara
Hertie School
Berlin, Germany
manohara
@hertie-school.org

Slava Jankin
University of Birmingham
Birmingham, UK
v.jankin
@bham.ac.uk

Abstract

This paper presents a multilingual aligned corpus of political debates from the United Nations (UN) General Assembly sessions between 1978 and 2021, which covers five of the six official UN languages: Arabic, Chinese, English, French, Russian, and Spanish. We explain the preprocessing steps we applied to the corpus. We align the sentences by using word vectors to numerically represent the meaning of each sentence and then calculating the Euclidean distance between them. To validate our alignment methods, we conducted an evaluation study with crowd-sourced human annotators using Scale AI, an online platform for data labelling. The final dataset consists of around 300,000 aligned sentences for En-Es, En-Fr, En-Zh and En-Ru. It is publicly available for download.

1 Introduction

Multilingual corpora are valuable resources for natural language processing (NLP) research and applications, as they enable the development and evaluation of cross-lingual and low-resource models and systems. However, creating and maintaining large-scale and high-quality multilingual corpora is a challenging task, as it involves collecting, processing, and aligning texts from multiple languages and domains, while ensuring their accuracy, consistency, and relevance. In this paper, we align and evaluate a multilingual corpus that is based on the plenary sessions of the United Nations (UN) General Assembly, which is the main

organ of the UN where all member states have equal representation and voice. The plenary sessions are held every year and are translated and transcribed in the six official languages of the UN: Arabic, Chinese, English, French, Russian, and Spanish. These sessions cover a wide range of global issues, such as peace and security, human rights, development, climate change, and health, and reflect the views and positions of different countries and regions on these issues. Therefore, our corpus provides a rich source of multilingual texts within the political domain that can be used for various NLP tasks, such as machine translation, in-domain text classification, question-answering, and multilingual argument mining.

We describe the methods we used to collect, clean, segment, and align the plenary sessions across languages. We then cover our approach to bilingual sentence alignment using embeddings and Euclidean distance, and the special considerations and difficulties we encountered for the different languages. For example, we faced some challenges in aligning Arabic with the other languages, due to technical issues in converting the Arabic documents into a suitable format for alignment. We also noticed some differences in the order and structure of sentences across languages, which made the alignment more difficult. We discuss how we addressed these challenges and how we validated the quality of our alignment.

2 Related Work

Previous multilingual parallel corpora have been based on United Nations data, including the MultiUN (Eisele and Chen, 2010) and the United Nations Parallel Corpus (Ziems et al., 2016).

In the MultiUN Corpus, data was retrieved from the United Nations Official Document Sys-

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

tem (ODS), a web-based repository of official documents of the United Nations. The data was filtered by publication symbols, which are unique identifiers that indicate the issuing body, the type of document and the year of publication. The paper selected documents with publication symbols that correspond to official records and other parliamentary documents of the UN. The multilingual sentence alignment starts with pairwise alignments based primarily on sentence lengths and then on a dictionary. Pairwise alignments are available, and later merged into multilingual alignments across all six languages. The updated version of MultiUN, v2, contains documents up to and including 2011 (Chen and Eisele, 2012).

The United Nations Parallel Corpus (UNPC) is composed of official records and other parliamentary documents of the United Nations that are in the public domain. It contains sentence-level alignments for content between 1990 and 2014. The corpus contains 799,276 documents in six languages and contains 86,307 documents that have translations across all six languages. The sentence-level alignments were generated using GIZA++ (Och and Ney, 2003) to align sentences based on word co-occurrences. The authors validate their dataset through a number of machine translation baselines, with BLEU scores results varying between 29 and 61, depending on the language pair.

While robust, neither of these corpora offer a full evaluation of the accuracy and precision of their alignment, nor are they recent enough to include the later documents. Furthermore, sentence alignment methods have come a long way since their collection. More modern methods of multilingual sentence alignments are based on multilingual pretrained language models, such as mBERT, that can learn cross-lingual representations of sentences. These methods have been shown to outperform GIZA++ (Schwenk, 2018; Guo et al., 2018). Artetxe and Schwenk (2019) use a sequence-to-sequence architecture to train a multilingual sentence encoder on an initial parallel corpus. The encoder maps sentences from different languages into a shared embedding space, where similar sentences are close to each other. The authors then use a margin-based scoring method to measure the similarity between sentence embeddings. The authors evaluate their their method on three tasks, the BUCC mining task, the UN reconstruction task, and the ParaCrawl filtering task, and show that the

proposed method outperforms existing methods on all three tasks by a large margin.

3 Corpus Collection

The United Nations (UN) plenary meetings are meticulously recorded in each of the six official UN languages, making them an ideal source for a multilingual corpus. The records are then made available on the official website¹, in the form of PDF files, separated by language and session. All the documents are public domain. We downloaded the documents in all 6 languages and converted them using an OCR-based tool in a pdf editor, discarding pictures, tables and style markers. In total, we processed 2113 documents between 1978 and 2021.

However, due to the age of the documents and the limitations of the OCR-based tool, we were unable to convert enough Arabic-language documents for use in alignment. Furthermore, the documents we did manage to convert were of poor quality. As a result, we were unable to align the Arabic sentences and eventually discarded the Arabic language documents until such a time that we can properly convert them.

We then processed these documents using the Natural Language Toolkit (NLTK) package (Bird et al., 2009). We used the toolkit to separate the documents into individual sentences, as the documents only provide paragraph boundaries. We then tokenised the sentences and removed stop words to made them ready for alignment.

4 Sentence Alignment

The documents described in Section 3 are not translated at the sentence level, but rather at the level of individual speeches taken as a whole. This means that each speech in one language has a corresponding speech in another language, but not necessarily each sentence. Therefore, in order to create a parallel corpus at the sentence level, we need to match each sentence in one language with the equivalent sentence in another language. This is a challenging task, as the sentences may not have the same order, structure, or length across languages. Furthermore, translations are not always a one to one mapping. Sometimes a sentence can be represented by multiple sentences in the other language, or multiple sentences can be condensed

¹<https://gadebate.un.org/en>

into a single sentence in another. Therefore, a simple probabilistic model based on sentence length would fail across languages with different scripts and language families. To solve this problem, we use a semantic similarity approach that aligns the sentences based on their meaning and content, rather than their form or position. We do this by using word vectors to represent the meaning of each sentence as a numerical vector. Then, we calculate the euclidean distance between the vectors of each language pair. The sentence pair with the smallest distance is the correct match.

For word vectorisation, we used Language-Agnostic Sentence Representations (LASER)², an open-source NLP toolkit developed by Facebook and trained on the Tatoeba corpus³. LASER performs sequence-to-sequence processing with an encoder-decoder approach. The encoder network, which is used to generate the embeddings we need, is a five layered bi-directional Long-Short-Term Memory (BiLSTM) network whose input is a string and output is a fixed-size vector in a 1024 dimensional space. Crucially, this space is shared by all languages, meaning that sentences with similar meaning in two different languages would be mapped to very near points in the space, regardless of how different the languages are.

The vectors are normalised and stored in a matrix, where each row represents a sentence and each column represents a language. We calculate the Euclidean distance between each sentence in one language and a window of 25 sentences in another language for each language pair. We select the pair of sentences with the smallest distance as the match. The window size is implemented to decrease time complexity as well as improve accuracy by not considering sentences too far away to have been the intended translation. This is done to prevent long, vague sentences that may be close to several other sentences from being matched numerous times, while also allowing for genuine cases where a sentence in one language has legitimately been represented by multiple sentences in the other. We also perform anchoring, where we identify special entities such as dates and numbers, and include only sentences in the target language that contain the same terms to be considered for matching.

²<https://github.com/facebookresearch/LASER>

³<https://tatoeba.org/en/>

Table 1: Statistics of Pairwise Aligned Sentences

	Sentences	Source Tokens	Target Tokens
En-ES	322,379	8,051,597	8,782,297
En-FR	325,968	8,145,802	8,885,067
En-ZH	300,281	6,849,901	6,503,222
En-RU	316,031	7,938,417	6,849,994

5 Validation

After matching, we performed a simple validation by training a linear regression model to predict sentence length for a translation in a target language, based on the length of the original sentence in English. We use this model to estimate the likelihood that a target language sentence is the correct translation for an English sentence. If the other language sentence length is either less than 50% or more than 150% of how much it is predicted to be by the model, it is discarded. We keep the remaining matches and add them to the corpus.

The statistics for all validated language pairs are presented in Table 1, which shows the number of sentences for each language pair, along with the number of tokens for each of the language pairs.

6 Evaluation

We evaluated the quality of our final validated dataset using crowd-sourced human annotations. To obtain reliable and consistent evaluations, we used Scale AI⁴, an online platform whose purpose is to generate labelled datasets for training AI models. Scale AI allows for the labelling of data such as images, videos, texts and 3D models.

We uploaded our parallel documents to Scale AI and requested the annotators to mark the sentences that are translations of each other in each language pair. We also provided them with clear instructions and examples of how to perform the task. We received the annotations from Scale AI in a JSON format, which we converted into a tab-separated format for further analysis. Scale also selects a “training set” of 20 sentence pairs, which it chooses from the corpus, for its crowd-sourced users, and discards results from users that perform below a threshold of 70% on the training set.

We designed the task to present annotators with two sentences, the “original” sentence in English and the “target” sentence in one of the target languages: French, Russian, Spanish or Chinese. The

⁴<https://scale.com/>

annotators were asked to read the sentences carefully and decide whether or not the two sentences are a match, a partial match, no match, or if they were unsure. In the instructions, the annotators were given detailed guidelines about what constitutes a match. If two sentences are direct translations of each other, or if all the information in the target sentence is present in the original sentence, they are considered to be a match. Furthermore, if the target sentence conveys the full meaning of the original sentence, annotators are to consider them a match. Partial matches occur when some information in the target sentence is not present in the original sentence, or vice versa. If the sentences are neither a full match nor a partial match, then annotators were to choose “no match”. We also included an “unsure” option, and discarded any responses that included it. Figures 1 and shows an example of the instructions for the English–French evaluation, the way they appear next to each sentence pair in the task. In addition to these instructions, further workflow directions were made available.

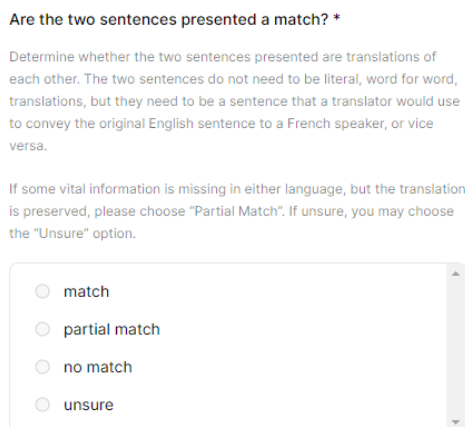


Figure 1: Instructions as they are presented to annotators alongside each sentence pair.

The task required 3 reviews per sentence, meaning three different annotators had to agree on a label for it to be accepted. Annotators were required to have a basic proficiency in the source language and native proficiency in the target language. However, due to crowd-sourcing, there was no way to verify their actual proficiency. The agreement between reviewers was pretty high, with Cohen’s Kappa at 0.87 across all four language pairs. Furthermore, the evaluators found that over 80% of the presented sentences were a match, and less 5% were completely unaligned.

The number of evaluated sentences varied across languages, as it depended on the number of available annotators that Scale was able to train for each task. As a result, while we only had 1000 sentences evaluated for English–French, we managed to evaluate upwards of 6000 sentences for English–Chinese. Table 2 shows the number of sentences we aligned, and the percentage of this total that we managed to evaluate.

Table 2: Percentage of Sentence Pairs Evaluated Across Languages

	Evaluation Set	Percentage of Total
En–Fr	1000	0.3%
En–Es	7000	2.3%
En–Zh	6000	2%
En–Ru	9000	3%

The evaluators found that on average, over 85% of aligned sentences were a complete match, with around 6% of sentences being completely misaligned. The English–Spanish language pair had the highest percentage of correctly aligned sentences, at 91.4% of sentences being a total match. Conversely, the English–Russian language pair showed the highest number of misalignment, with only 78.5% of sentences matching. Table 3 shows the percentage of correctly aligned sentences by language pair.

Table 3: Scale AI Evaluation of Alignment

	Complete Match	Partial Match	No Match
En–Fr	86.4%	8.5%	5.1%
En–Es	91.4%	2.6%	6%
En–Zh	86.3%	9%	4.7%
En–Ru	78.5%	13%	8.5%

7 Limitations

We acknowledge some limitations of our dataset that we aim to overcome in future iterations. One of the main limitations is that we could not include Arabic as one of the languages in our corpus, due to technical difficulties in converting the Arabic documents into a suitable format for alignment. This means that our dataset does not cover all six official languages of the United Nations, and thus misses an important and widely spoken language in the world. We hope to solve this problem by finding a more reliable way to process the Arabic documents and align them with the other languages. Another limitation of our dataset is

that we relied on crowd-sourcing for evaluating the quality of our alignment. While crowd-sourcing is a convenient and cost-effective way to obtain human judgments, it also comes with some drawbacks, such as inconsistency and bias among the annotators. We tried to mitigate this issue by auditing the results and filtering out the outliers, but we could only review a small fraction of the evaluations. Therefore, our evaluation may not reflect the true quality of our dataset, and may be influenced by the subjective opinions of the annotators. We plan to address this issue by conducting a more rigorous and systematic evaluation of our dataset, using multiple sources of human feedback and objective metrics.

8 Conclusion

In this paper, we introduced a novel parallel corpus that consists of texts from the plenary sessions of the United Nations General Assembly. Our corpus covers five languages: English, French, Spanish, Russian, and Chinese. We described the process of extracting and preprocessing the sentences from the original documents, and aligning them based on semantic similarity using a state-of-the-art cross-lingual sentence encoder. We evaluated the quality of our dataset using two methods: a simple validation that uses a regression model to predict sentence length based on the source language and the target language, and a crowd-source human evaluation that measures the accuracy and precision of our alignment.

The resulting aligned dataset has a high degree of accuracy across languages, and can be used for various natural language processing tasks, such as machine translation, cross-lingual information retrieval, and multilingual text summarisation.

Our work contributes to the field of multilingual natural language processing by providing a large-scale and high-quality parallel corpus that covers multiple languages in the field of political discourse and debate. We believe that our corpus can facilitate the development and evaluation of cross-lingual models and applications. In the future, we plan to solve the problem of Arabic-language documents that prevented us from completing our dataset for all six official languages of the United Nations. We also intend to extend our corpus to include more languages and more sources of multilingual texts. The current version of our dataset

is available for download⁵.

9 Acknowledgements

This project has received funding from the European Union’s Horizon Europe research and innovation programme under Grant Agreement No 101057131, Climate Action To Advance HeaLthY Societies in Europe (CATALYSE).

References

- Artetxe, Mikel and Holger Schwenk. 2019. Margin-based parallel corpus mining with multilingual sentence embeddings. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3197–3203, Florence, Italy, July. Association for Computational Linguistics.
- Bird, Steven, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. O’Reilly Media, Inc.
- Chen, Yu and Andreas Eisele. 2012. MultiUN v2: UN documents with multilingual alignments. In *International Conference on Language Resources and Evaluation*.
- Eisele, Andreas and Yu Chen. 2010. MultiUN: A multilingual corpus from united nation documents. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC’10)*, Valletta, Malta, May. European Language Resources Association (ELRA).
- Guo, Mandy, Qinlan Shen, Yinfei Yang, Heming Ge, Daniel Cer, Gustavo Hernandez Abrego, Keith Stevens, Noah Constant, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. Effective parallel corpus mining using bilingual sentence embeddings. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 165–176, Brussels, Belgium, October. Association for Computational Linguistics.
- Och, Franz Josef and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29, 03.
- Schwenk, Holger. 2018. Filtering and mining parallel data in a joint multilingual space. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 228–234, Melbourne, Australia, July. Association for Computational Linguistics.
- Ziemski, Michal, Marcin Junczys-Dowmunt, and Bruno Pouliquen. 2016. The united nations parallel corpus v1.0. In *International Conference on Language Resources and Evaluation*.

⁵https://github.com/KrishnaM313/UN_multilingual_corpora

Generating subject-matter expertise assessment questions with GPT-4: a medical translation use-case

Diana Silveira

University of Lisbon, Portugal
Unbabel, Lisbon Portugal
INESC-ID, Lisbon Portugal
dianasilveira
@edu.ulisboa.pt

Marina Sánchez-Torrón

Unbabel
marina.sanchez@
unbabel.com

Helena Moniz

University of Lisbon, Portugal
INESC-ID, Lisbon Portugal
helena.moniz@
edu.ulisboa.pt

Abstract

This paper examines the suitability of a large language model (LLM), GPT-4, for generating multiple-choice questions (MCQs) aimed at assessing subject-matter expertise (SME) in the domain of medical translation. The main objective of these questions is to model the skills of potential subject-matter experts in a human-in-the-loop machine translation (MT) flow, to ensure that tasks are matched to the individuals with the right skill profile. The investigation was conducted at Unbabel, an artificial intelligence-powered human translation platform. Two medical translation experts evaluated the GPT-4-generated questions and answers, one focusing on English–European Portuguese, and the other on English–German. We present a methodology for creating prompts to elicit high-quality GPT-4 outputs for this use case, as well as for designing evaluation scorecards for human review of such output. Our findings suggest that GPT-4 has the potential to generate suitable items for subject-matter expertise tests, providing a more efficient approach compared to relying solely on humans. Furthermore, we propose recommendations for future research to build on our approach and refine the quality of the outputs generated by LLMs.

1 Introduction

This work presents an approach for developing an assessment tool to evaluate the subject-matter expertise (SME) of professional translators in the field of medical translation, using the large language model GPT-4. As MT becomes more predominant in translation scenarios, including specialized fields, the need for skilled experts who can identify and address quality concerns in MT-generated output is proportionately increasing.

Specialized translators require SME, which involves extensive knowledge and proficiency in a specialized domain in both the source and target languages of the relevant language pair. To measure and evaluate SME effectively, a high-quality test for translators should incorporate translation questions that assess the translator’s proficiency in the target language, as well as questions that evaluate their domain-specific expertise in the source language from which they are translating (Montalt, 2007).

Evaluating the SME of potential experts in the loop across different language pairs and domains poses challenges. Implementing a system that uses SME tests to pre-screen and match subject-matter experts with MT texts in the same subject matter could improve quality in a human-in-the-loop flow. However, when done entirely by humans, creating and maintaining a comprehensive and up-to-date question bank for a wide array of language pairs and domains can become expensive and time-consuming. Other challenges are listed in Section 5.

This paper proposes a methodology for automating the creation of SME tests using GPT-4, focusing on the field of medical translation, and the English–European Portuguese and English–German

language pairs. Opting for a multiple-choice questionnaire format allowed for the automation of both test generation and test grading, which increases the speed and scalability of the assessment process, while decreasing the costs. Apart from possessing high levels of objectivity, MCQ tests facilitate a fast and effective assessment process, while also being able to cover a broad range of topics (Jovanovska, 2018). Our main objective with the question banks is to distinguish between experts and non-experts, not to evaluate levels of expertise among the expert test-takers, although that might constitute a future point of research (See Section 5 for insights on future work).

We conducted a quality assessment of the generated questions and sets of answers, using an evaluation scorecard structured around five evaluation criteria (detailed in Section 2), which was completed by medical expert translators. Our analysis of the results provides insights into GPT-4’s capability in generating high-quality medical expertise test items for English–European Portuguese and English–German. All the data related to this work, including question banks, evaluation and prompts, is available on GitHub.¹

2 Study setup

Our plan for implementing SME tests as part of Unbabel’s framework for matching human experts to translation tasks involves generating different assessments for each test-taker, by randomly choosing a set number of questions from a large question bank. As such, our approach centered on compiling a large question bank generated by GPT-4, rather than individual tests, for each language pair. For each language pair, we invited a medical expert translator with more than ten years of professional experience to evaluate the generated questions and answers. The evaluators had no connection to the study and were paid according to their hourly rates with no time restrictions to conduct the experiment.

2.1 Question bank typology

For each language pair, we generated four separate question banks, of 50 unique questions each. Each question bank was generated with a specific prompt. The question banks are categorized by:

For each language pair, we generated four separate question banks, of 50 unique questions each. Each question bank was generated with a specific prompt. The question banks are categorized by:

1. Topic: each focuses on a different area/type of document within the medical translation field. These are: a) clinical trials and clinical trial protocols; b) general medical information; c) clinical studies and d) terminology translation (which encompasses the previous three topics).
2. Language: they are either a source language only question bank or a translation question bank.
3. Question type: each bank is based on a different format of MCQ: a) multiple choice with four options, one of which is correct, b) alternate-choice questions/true or false, and c) fill-in-the-blanks questions, with four options.

Each medical translation expert, one per language pair, assessed all four question banks. Three of those question banks —fully in English, the source language— were shared across both language pairs. Only one question bank included translation-related questions.

2.2 Prompts and model parameters

The prompts were fed onto the GPT-4 model on November of 2023 via OpenAI’s Playground, with specific parameters to shape the output: temperature, which adjusts the randomness of the model’s predictions, was set to 0.4 to enhance the accuracy of the response; Top P was set to 1, ensuring the model’s predictions included the whole range of possibilities; and presence penalty, which discourages repetition, was set to 0.5, encouraging the model to introduce new ideas and topics for a varied response.

The process of constructing the ideal prompt was incremental, performed in a trial and error manner. We started with a simple instruction: “*I want to test the subject matter expertise of translators in the domain of [chosen domain]. Create a questionnaire containing [chosen number of items] multiple choice questions.*” On subsequent iterations, several instructions were added to the prompts, in order to curtail issues as they arose. To elevate the difficulty of the question banks, we

¹<https://github.com/mstorrton/subject-matter-expertise-assessment-questions-with-GPT-4>.

	Question type	Subfield	Language
QB1	Multiple choice A), B), C), D)	Clinical trials and clinical trial protocols	Source language (EN)
QB2	Alternate choice	General medical information	Source language (EN)
QB3	Fill in the blanks	Clinical studies	Source language (EN)
QB4	Multiple choice, true and false, fill-in-the-blanks and alternate choice	Terminology translation	EN–PT
QB5	Multiple choice, true and false, fill-in-the-blanks and alternate choice	Terminology translation	EN–DE

Figure 1: Table 1. Question banks divided by categories.

instructed the model to produce items with a level of complexity that would make it difficult for non-experts to answer the questionnaire correctly, and to clearly identify the correct answer in the choices given.

To ensure the output matched our language and format expectations, we specified, for QB1, QB2 and QB3, that test items were to be fully in English, with no translation items, and detailed the format of the stem and answer choices, according to each QB’s typology (See Table 1).

In order to increase the quality of the distractors, we added the instruction to include plausible distractors, which we further detailed as: “*the answer choices should be similar to each other and the question in category, morphology or syntax*”. Additionally, we instructed the model to provide distractors with the same length and complexity, as well as to ensure that the correct answer and the question stem did not share words with the same word root.

For generating QB4 and QB5, we developed glossaries with domain-specific word pairs using GPT-4 beforehand, to ensure the relevance of the terms used in these question banks. Each glossary contained 50 word pairs related to general medical information, clinical trials and clinical studies,

with each word pair used to create one test item.

We used zero-shot prompts for all the question banks, except for QB4, the English–European Portuguese translation question bank, for which we used a few-shot prompt, containing two examples of the ideal type of output. We used few-shot prompting only on QB4, as a means of comparing the quality of the distractors compared to zero-shot prompting; few-shot prompting yielded better distractors (See Section 5.5). As GPT-4 is a commercial model, with charges based on the combined number of input and output tokens, prioritizing zero-shot prompts is generally more cost-effective. To see the complete prompts and respective outputs, refer to the GitHub link provided in the Introduction.

2.3 Evaluation scorecard

We created a question scoring system based on five multiple choice quality criteria, which we curated based on the works of Town (2014) and Jovanovska (2018).

1. Question accuracy: is the question worded clearly and unambiguously, so that the correct answer could be clearly identified by an expert?
2. Correct answer factuality: is the correct answer choice, also known as *key*, scientifically true?
3. Non-ambiguous answer choices: is there more than one correct answer?
4. Prevalence of correct answer: is the correct answer the most prevalent and commonly applied option in the context of the question?
5. Plausible distractors in the answer choices: do the incorrect answer choices constitute plausible distractors for non-expert test-takers?

The scoring system is binary, relying on “Yes” or “No” responses to evaluate each question and its answer choices according to the above criteria. The scorecard was created using Google Sheets, containing one test item per row and one question bank per sheet. The subject-matter experts had access to the test items in the following format: question, answer choices, key (selected by the model); followed by the five criteria presented above and ending with a column for comments and a column for the score. They assess and evaluate each question and set of answer choices on each criterion using a drop-down menu with “Yes” and “No” options, which results in an automatic score based on their evaluation. To further clarify, for criterion 3

“Is there more than one correct answer?”, the ideal answer would be “No”, as ambiguity is not desired in these types of tests. In the case of criterion 5, “Do the incorrect answer choices constitute plausible distractors for non-expert test-takers?”, the ideal answer would be “Yes”, as this would prevent test takers from achieving high scores simply from “guessing”.

3 Results

After generating the question banks, the next step was assessing and scoring their quality. Section 3.1 presents the overall score, results and considerations for English–European Portuguese, conducted by Evaluator A, while Section 3.2 does the same for English–German, conducted by Evaluator B.

3.1 English–European Portuguese question banks scores

On Table 2, it can be observed that “Question accuracy” achieved the highest possible score on every question bank. Conversely, “Plausible distractors in the answer choices”, is the overall lowest scoring parameter across all four question banks. When this was the case, it was often because the correct answer would be a term that shared the same root as a word in the question, but the distractors did not, as you can see in the following example:

What is the medical term for inflammation of the pancreas?

- A) Pancreatitis
- B) Gastritis

From this, we infer that GPT-4 frequently fails to follow the instruction “Do not include correct answers that share the same root as words in the stem”, when given a zero-shot prompt. However, when the model was given a few-shot prompt, as is the case with QB4, the plausible distractor category achieved the highest score. Despite the plausibility of the distractors, Evaluator A stated that QB4 had a few items with ambiguous answer choices, meaning that more than one answer choice could be considered correct. On QB1, QB2 and QB3, “Correct answer factuality” scored highly, and so did “Prevalence of correct answer”: in the majority of test items, the answer indicated by the model as the key (correct answer) was the

EN-PT Medical translation question banks	QB1	QB2	QB3	QB4
Question accuracy	100	100	100	100
Correct answer factuality	100	100	96	80
Non-ambiguous answer choices	96	100	98	74
Prevalence of correct answer	98	100	98	88
Plausible distractors	96	80	80	100
Average	98	96	94.4	88.4

Figure 2: Table 2: Scores of the EN-PT question banks. QB1 - Clinical trials and clinical trial protocols QB2 – General medical information QB3 – Clinical studies QB4 – EN-PT medical terminology

most prevalent within the context of the question, according to Evaluator A.

3.2 English–German question banks scores

Similarly to English–European Portuguese, “Question accuracy” achieves a perfect score on all four question banks for English–German. “Correct answer factuality” scored highly across the question banks, with a few exceptions on QB1 and QB3. These two question banks also presented the highest amount of ambivalent answer choices. When there were more than two possible correct answer choices, the answer identified as correct was still the most prevalent option in QB1, QB2 and QB4, which demonstrates high precision by GPT-4 when identifying factual answers, but a lower capacity for providing answer choices that are both plausible distractors and unambiguously incorrect. This can be seen in the following example:

The German translation for “Vertigo” is:

- A) Vertigo

- B) Schwindel
- C) Schwindelgefühl
- D) Vertigokrankheit

Option C) was selected as the key by GPT-4, but options A) and B) were also correct translations. Once more, “Plausible distractors” was the lowest scoring parameter across all question banks, with lower scores on QB2 and QB3. Still, it is worth noting that its score, across all question banks, never reaches below 78 points out of 100.

EN-DE Medical translation question banks	QB1	QB2	QB3	QB4
Question accuracy	100	100	100	100
Correct answer factuality	92	100	98	98
Non-ambiguous answer choices	96	98	98	86
Prevalence of correct answer	100	100	98	100
Plausible distractors	92	80	78	92
Average	96	95.6	94.4	95.2

Figure 3: Table 3: Scores of the EN-DE question banks. QB1 - Clinical trials and clinical trial protocols
QB2 – General medical information
QB3 – Clinical studies
QB4 – EN-DE medical terminology

3.3 GPT-4 generated glossaries

As mentioned in Section 2.2, the prompts for the translation question banks included glossaries of 50 word pairs, one for each test item. The English-German glossary was generated by GPT-4 in the following way: we requested a set of 50 medical terms, in English, related to the topics mentioned in Section 2.1. From that output, we asked the model to replace repetitive or irrele-

vant items, which we singled out from the original list, using the instruction “Replace items [number of each item in the list]” Finally, we asked the model to translate the medical terms into German, resulting in a curated glossary of 50 word pairs. The evaluator considered all the translations in the glossary correct, but pointed out two repeated items that had not been previously detected. The English-European Portuguese glossary was only partly generated by GPT-4: 26 word pairs were obtained from a publicly available glossary on the medical subfield of clinical trials, L10N Studio. The remaining 24 word pairs were generated in the same way described above. This division showed clear results: the word pairs generated by GPT-4 showed some terminology and mistranslation issues. For example, “clinical pharmacology study” was translated as “estudo clínico de medicamento”, when it should be “estudo clínico de farmacologia”, and “particle therapy” was translated as “terapêutica de partículas” when it should be “terapia de partículas”. The word pairs extracted from the verified source, on the other hand, were deemed much more accurate by our SME evaluator, with only one out of the 26 term pairs not considered the ideal translation.

3.4 Overall results

Question bank overall quality scores	
QB1 – Clinical trials and clinical trial protocols	97
QB2 – General medical information	95.8
QB3 – Clinical studies	94.4
QB4 – EN-PT medical terminology	88.4
QB5 – EN-DE medical terminology	95.2

Figure 4: Table 3: Overall scores of the five GPT-4 generated question banks. The maximum possible score for each question bank is 100.

Table 3 shows the overall quality of each of the five question banks generated by GPT-4 for evaluating medical translators’ subject matter expertise. (Note: for the score of QB1, QB2 and QB3, we calculated the average of the scores of both evaluators, when they differed). English-European Portuguese (four question banks combined) has an overall score of 94.2%, while English-German has an overall score of 95.3%.

The results reflect the high level of quality and practical applicability of the generated question banks. In terms of perceived level of difficulty, the evaluators gave the question banks an average of 3.5 out of 5 (1 being very easy and 5 very difficult). For more on the difficulty dimension, see Section 5.4.

4 Limitations

Despite its preliminary positive findings, this study presents several limitations. Firstly, the topics chosen for the question banks only represent a very small portion of medical knowledge, and only address a few of the most commonly translated medical documents. As mentioned in Section 4, reproducing the study with different and perhaps less common language-pair combinations is likely to produce different results. We hypothesize that the less common the language-pair combination, the lower the quality achieved by GPT-4. Dac Lai et al (2023) state that there is a decrease in performance for languages other than English in natural language processing (NLP) tasks, which might be verified in the use-case of MCQ automation. The same can be said for less common language varieties. Additionally, the evaluation of the GPT-4 output was done by only one expert per language pair. The sample size evaluated by each expert was substantial (200 question stems and 200 sets of answer choices), but extending the evaluation process to more experts can strengthen the validity of the results. Furthermore, the ontology of subject matter expertise is vastly complex and multifaceted (Collins and Evans, 2007; Shavelson, 2010) and this paper does not intend to claim that the measurement of subject matter expertise can be fully judged by the results of MCQ assessments. The resulting MCQs of this study are tailored for specific use-cases, not to measure competency in general. They are also designed to be part of a larger assessment process, in which other specific tasks (such as reviewing a specialized machine translated text, for instance) contribute to a more accurate representation of the expertise level of the human-in-the-loop.

5 Conclusion and future work

MCQ automation using large language models has been a prevalent topic of research in a wide range of fields, such as reading comprehension (Sayin

et al. 2024), vocabulary testing (Wang et al. 2024), programming (Doughty et al. 2024) and medical education (Kiyak, 2023), among others. In this study, we observed promising results regarding GPT-4's capability to generate SME tests for specialized translators, in the medical domain, with the English–European Portuguese and English–German language pairs. In order to verify the applicability of the findings in this study, it is recommended to replicate the study with other language pairs and subject-matter domains. While we selected GPT-4 to perform the study, the same methodology might yield high-quality results with other LLMs. The objective of this study was to determine the viability of automating the generation of MCQs for assessing and labeling the skills of expert translators at Unbabel, to match them to specialized tasks requiring those skills. The overall quality of the four question banks combined was 94,2% for English–European Portuguese and 95.3% for English–German. This indicates that including GPT-4-generated MCQs in our expertise assessment process is a viable option. Our initial aim for the generated questions is to differentiate non-experts from experts. In the future, it may be interesting to assign different levels of expertise based on the percentage of correct answers. To achieve this, we might need to introduce more challenging questions and distractors. This will be considered once we analyze the difficulty of the current question banks.

What follows are recommendations for improving the relevance of the output from an LLM for the use case presented in this paper: generating SME test items to be ultimately used as an optimization method of task assignment in expert-in-the-loop translation flows. They can also be adapted for other contexts, MT related or otherwise. These SME test items might be relevant, for instance, in traditional translation workflows, research surveys to gather data, or in businesses, in the context of assessments and job interviews. Another relevant use case is the integration of the generated test items in self-directed learning methodologies (Loeng, 2020), either in corporate contexts or within freelance translator and reviewer training.

5.1 Language-pair glossaries

When creating the prompt for QB4 and QB5 (the translation question banks), instead of utilizing

glossaries generated by GPT-4, we recommend curating an up-to-date glossary of terms taken from one or more reliable and accredited sources; the number of word pairs to include should be the same as the number of questions requested in the prompt. Requesting a mixed format question bank (including different typologies of multiple-choice questions) resulted in varied and diversified question banks, as was the case for the translation question banks. Including high-quality specialized glossaries does constitute an extra step before crafting the prompts, but it guarantees a superior result and less intervention when it comes to the human step of reviewing, validating and (potentially) correcting the question banks.

5.2 Inclusion of relevant word pairs

For entities and organizations that have similar types of specialized documents and materials with which they work regularly, we recommend extracting the most common and relevant domain-specific terms found in the translated content, and adding them to the word-pair glossaries. That way, the tests generated by GPT-4 or other LLMs will become more tailored to the organization's workflow.

5.3 Elimination of alternate-choice question format

It is less likely to have plausible distractors with only two answer choices (Towns, 2014) and we verified that the distractors provided by GPT-4 often decreased the level of difficulty, making it easier for non-experts to guess the correct answer. Instead, we consider it is more beneficial to replace the alternate-choice question bank, QB2, with a classic four-option MCQ question bank, maintaining, however, the same topic (general medical information). On QB4 and QB5 (the translation section), we would likewise remove the instruction for including alternate-choice questions.

5.4 Evaluation of item difficulty and discrimination

A viable next step for this research would be to evaluate the level of difficulty of the generated test items. For this, we would employ a difficulty index, which would indicate the percentage of test-takers who answered each question correctly, as well as a discrimination index, which calculates the relationship between each individual test taker's test item score with the overall scores

of all test takers, allowing each test item to discriminate between high and low scorers (Hingorjo and Jaleel, 2012). With this knowledge, we would be able to establish a pass/fail threshold with a percentage (to be defined) of correctly answered test items that distinguishes an expert test-taker from a non-expert, when it comes to the experts-in-the-loop who would be assigned to Unbabel's domain specific translation tasks. With the difficulty index, we could also determine which test items prove to be extremely easy or extremely difficult (therefore not good indicators of SME) and potentially eliminate them from our pool of questions. This is a necessary step before implementing the tests at Unbabel and measuring their impact on the company's workflows, to ensure that we are using only the most accurate and appropriate test materials.

5.5 Few-shot prompting

We strongly recommend using few-shot prompts containing at least two examples of test items with plausible distractors. For this study, we ascertained that distractors considered plausible had to be either semantically plausible (usually in the same category as the key) or morphologically plausible, which means they would also contain terms sharing the same word-root as the key or the question stem. This was most successfully achieved with the use of few-shot prompting on QB4, which led to the conclusion that this prompting technique is the most adequate to generate high quality distractors for the present use case of question bank automated generation.

5.6 LLM and language-pair diversity

Finally, this study should be replicated in the future with different LLMs and language pairs, as well as different areas of specialized translation, to extend its findings and further assess the validity of this type of methodology.

Acknowledgements

This work was developed within the scope of the project n° 62 - "Center for Responsible AI", financed by European Funds, namely "Recovery and Resilience Plan - Component 5: Agendas Mobilizadoras para a Inovação Empresarial", included in the NextGenerationEU funding program and was partially founded by FCT, Fundação para a Ciência e a Tecnologia, under project UIDB/50021/2020

References

Doughty, Jacob, Zipiao Wan, Anishka Bompelli, Jubahed Qayum, Taozhi Wang, Juran Zhang, Yujia Zheng, Aidan Doyle, Pragnya Sridhar, Arav Agarwal, Christopher Bogart, Eric Keylor, Can Kultur, Jaromir Savelka, Majd Sakr. 2024. A comparative study of AI-generated (GPT-4) and human-crafted MCQs in Programming Education. *Proceedings of the Australian Computing Education Conference 2024*.

Hingorjo, Mozaffer Rahim, and Farhan Jaleel. 2012. Analysis of One-Best MCQs: the Difficulty Index, Discrimination Index, and Distractor Efficiency. *Journal of the Pakistan Medical Association*. 62(2): 142–7.

Jovanovska, Jasmina. 2018. Designing effective multiple-choice questions for assessing learning outcomes. *Infotheca - Journal for Digital Humanities* 18(1): 25–42.

L10N Global. 2018. 30 Common clinical trial terms in English and Portuguese (PT and BR). Accessed on 9/10/2023. <https://www.l10nglobal.com/en/news/30-common-clinical-trial-terms-english-and-portuguese-pt-and-br>

Loeng, Svein. 2020. Self-directed learning: A core concept in adult education. *Education Research International*, 2020: 1-12.

Lai, Viet Dac, Nghia Trung Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Derroncourt, Trung Bui, Thien Huu Nguyen. 2023. ChatGPT Beyond English: Towards a Comprehensive Evaluation of Large Language Models in Multilingual Learning. *Paper presented at the Conference on Empirical Methods in Natural Language Processing 2023*.

Montalt, Vicent, and Maria González-Davies. 2006. *Medical Translation Step by Step: Learning by Drafting*. Oxfordshire, England: Routledge.

OpenAI. Accessed in November of 2023. "Playground." <https://playground.openai.com>.

Sayin, Ayfer, et al. 2024. Using OpenAI GPT to Generate Reading Comprehension Items. *Educational Measurement Issues and Practise*. 43(1): 5–18.

Shavelson, Richard J. 2010. On the measure-

ment of competency. *Empirical Research in Vocational Education and Training*. 2(1): 41–63.

Wang, Qiao, Ralph Rose, Naho Orita, and Ayaka Sugawara. 2024. Automated Generation of Multiple-Choice Cloze Questions for Assessing English Vocabulary Using GPT-turbo 3.5. *Proceedings of the Joint 3rd International Conference on NLP4DH and 8th IWCLU*.

Prompting Large Language Models with Human Error Markings for Self-Correcting Machine Translation

Nathaniel Berger^{a*} and Stefan Riezler^{ab}

Computational Linguistics^a & IWR^b
Heidelberg University
69120 Heidelberg, Germany
berger@cl.uni-heidelberg.de
riezler@cl.uni-heidelberg.de

Miriam Exel^c and Matthias Huck^c

SAP SE^c
Dietmar-Hopp-Allee 16
69190 Walldorf, Germany
miriam.exel@sap.com
matthias.huck@sap.com

Abstract

While large language models (LLMs) pre-trained on massive amounts of unpaired language data have reached the state-of-the-art in machine translation (MT) of general domain texts, post-editing (PE) is still required to correct errors and to enhance term translation quality in specialized domains. In this paper we present a pilot study of enhancing translation memories (TM) produced by PE (source segments, machine translations, and reference translations, henceforth called PE-TM) for the needs of correct and consistent term translation in technical domains.

We investigate a light-weight two-step scenario where, at inference time, a human translator marks errors in the first translation step, and in a second step a few similar examples are extracted from the PE-TM to prompt an LLM. Our experiment shows that the additional effort of augmenting translations with human error markings guides the LLM to focus on a correction of the marked errors, yielding consistent improvements over automatic PE (APE) and MT from scratch.

1 Introduction

Technical translation at large enterprises involves a large number of translation domains, for which translation memories and terminologies need to

be maintained to support multi-domain MT systems and human post-editors in producing contextually adequate and consistent translation of technical terms (Exel et al., 2020). In this paper, we ask if ongoing human post-editing efforts that produce large databases consisting of source segments, machine translations, and reference translations, can be enhanced by light-weight human error markings. This could then be used to teach a translation system a focused self-correction of marked erroneous tokens from similar examples with error markings and corrections found in the PE-TM. Such a setup could complement translation memories and terminology databases by up-to-date and domain-specific information in the PE-TM, and be used in a scenario where a user marks errors in MT hypotheses. In-context examples with high source-side similarity are then extracted from the PE-TM to prompt an LLM to focus on a correction of the marked error interactively.

We present a pilot study where we construct a PE-TM for the IT domain, which is augmented by human error markings on machine translations. While for training purposes, error markings for the PE-TM could be obtained by automatic matching against human post-edits, this cannot be done at test time. We envisage a scenario where the error markings in the PE-TM are obtained by direct human annotation, simulating a realistic setup where a user only performs the light-weight task of error marking at test time. Such a scenario could be feedback collection in the publishing of raw-MT. Raw-MT could be shown to end-users who, if they notice an error in the translation, proceed to annotate tokens in the translation they perceive to be incorrect. The translation would then be flagged for review by a human translator, who then post-edits the translation and publishes their correction. This process results

*The work was done as part of an SAP sponsored PhD project of the first author.

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

Read the English text and the German translation hypothesis and then correct the output. Incorrect words are inside of tags '<bad> </bad>'. Please use this feedback in your correction. If the hypothesis is already correct, do not make any changes.

English: This environment variable can also be used to make sure that other operations are working on uploaded files, as well.

Hypothesis: Dieses <bad> Umweltvariable </bad> kann auch verwendet werden , um sicherzustellen , dass andere <bad> Operationen auf </bad> hochgeladene Dateien <bad> arbeiten </bad> .

German: Mittels dieser Umgebungsvariable kann auch sichergestellt werden, dass auch andere Operationen an hochgeladenen Dateien arbeiten können.

English: Some important environment variables used by KDE

Hypothesis: Einige wichtige <bad> Umweltvariablen </bad> , die von KDE verwendet werden

German: **Einige wichtige Umgebungsvariablen, die von KDE verwendet werden**

Figure 1: Example of a 1-shot prompt for English-to-German Translation. Error markings are inside bold faced tags <bad> </bad>. The demonstration example consists of a source segment in English (in green), a translation hypothesis in German (in blue), and a correction (in red). The test example shows a correction of the translation of "environment variable" from "Umweltvariable" into "Umgebungsvariable" learned by the LLM (in bold-faced red).

in the creation of (source, hypothesis, post-edit) triples with annotations for the PE-TM. This PE-TM is used to provide in-context examples for LLM correction of annotated translation hypotheses. For example, the end-user who annotates raw-MT could then immediately be shown a new translation that takes the error markings into account.

The results of our study show that selecting in-context examples based on similarity of source-side embeddings and providing error markings on hypotheses lets the LLM infer focused corrections of marked errors. Furthermore, overall translation quality is improved over few-shot prompt-based translation and over automatic post-editing. An example 1-shot prompt and error-marked output is given in Figure 1.

2 Prior Work

The last year has seen a progression of the translation capabilities of decoder-only LLMs, pre-trained on unpaired language data, from lagging behind supervised systems (Vilar et al., 2023) to matching their translation quality (Garcia et al., 2023), with only 5 examples of high-quality translation data used for in-context learning. However, MT in specialized domains still requires translation post-editing in order to correct errors and to enhance term translation quality. Raunak et al., (2023) recently showed that very large LLMs (OpenAI,

2023) can perform zero-shot automatic translation post-editing for general domain data, at the price of hallucinated edits. This makes this setup impractical if high precision in domain-specific translation is key. For these purposes, manually crafted glossaries (Vidal et al., 2022), dictionaries extracted in a separate step of unsupervised word-alignment (Ghazvininejad et al., 2023), or translation memories accessed with fuzzy matching (Moslem et al., 2023; Hoang et al., 2023), have been used to aid prompt-based MT. Our approach combines PE-TMs with light-weight human error markings, achieving improvements over both APE and MT from scratch.

The standard paradigm to incorporate token-level human error markings as learning signal is an adaptation of supervised learning from post-edits (see, for example, Turchi et al. (2017)) by penalizing erroneous tokens and rewarding correct tokens in a weighted maximum-likelihood objective (Marie and Max, 2015; Domingo et al., 2017; Petrushkov et al., 2018; Lam et al., 2019; Kreutzer et al., 2020; Berger et al., 2023). Most approaches are conceptualized as fine-tuning applications, with error markings obtained by automatic matching against human post-edits or by direct human annotation. The approach that is closest to our work is QuickEdit (Grangier and Auli, 2018). They train a model with separate encoders for source and error-marked

hypothesis in order to improve upon the initial hypothesis by avoiding the marked tokens. Similar to our approach, QuickEdit requires error-markings at inference time. While QuickEdit relies on supervised learning, our approach succeeds in teaching an LLM to avoid marked tokens from a few demonstration examples of similar error patterns.

More recent work by Xu et al., (2023) successfully uses feedback in form of error type and location that is predicted by a learned error pinpoint model. Their work focuses on general domain translation and quality-estimation type feedback, in difference to the focused error markings on technical terms that we are interested in. We plan an extension of our work in the direction of using learned error markings in future work.

Our work is furthermore related to the more general issue of self-correction capabilities of LLMs. Similar to the findings of Huang et al., (2023), our work shows that in order to qualify as a correction rather than a mere change, automatic self-correction in LLMs (Madaan et al., 2023; Pan et al., 2023) needs to be guided by an oracle. In our case, the oracle consists of feedback on the errors in translation outputs of an LLM, combined with a few examples of similar errors and their reference translations.

3 Data and Models

We collected English and German parallel data from open source software documentation and localization available on OPUS (Tiedemann, 2012), as this data comes closest to our domain of interest. We concatenated data from *GNOME*, *KDE4*, *KDEdoc*, *PHP*, and *Ubuntu* to create our data set and filtered them with the following methods: we removed those segments containing fewer than five words or more than 25, those identified by fast-Text (Joulin et al., 2017) as the wrong language, those with more than 20% of characters being non-alphanumeric, and those containing personally identifiable information. Of the remaining data, we selected a subset of 1,500 examples.

For the purposes of this experiment, we were interested in models that support prompt-based interaction. Furthermore, we are interested in the scenario where users use their judgment to guide a model towards a better translation based upon its original translation. Large language models lend themselves well to this interaction because the same model can be used with prompt-based interaction to produce the original translations as well as for pro-

viding extra information to aid in correction. These considerations decide in favor of using an LLM over a traditional encoder-decoder based model typically used in production scenarios. Therefore, we examine if the model that produced the hypothesis, Llama 13B (Touvron et al., 2023), can leverage feedback to correct its own mistakes. In addition to choosing this model because it supports prompting, it runs on a single GPU¹, and the model will remain available in the future for reproducibility. The Llama model was converted to Huggingface Transformers (Wolf et al., 2020) format for inference. We found that Llama 13B frequently copies hypotheses including the error tags to its output as it was instructed not to make changes if the hypothesis is acceptable as-is. We therefore post-processed outputs by removing tags by regex.

Additionally, we test GPT-3.5² in order to test a model larger than we can locally run. We use the "ChatCompletion" API, send the query as the 'user' message, and set the temperature to 0. All hypotheses were generated with the above models using greedy decoding.

4 Feedback Collection

In order to simulate the previously proposed scenario of an end-user who annotates raw-MT errors, we turn to paid annotators. We generate translations of English source sentences and provide them with only the source and the hypothesis, as would be the case when getting feedback on raw-MT. These are then paired with references to create our PE-TM.

4.1 Human Annotation

We hired three professional translators with expertise in the IT domain as annotators to provide token-level feedback on the translation hypotheses. Token-level feedback consisted of per-token binary quality judgements, OK/BAD. Annotators were provided English source sentences and German hypothesis translations in a custom annotation interface. Each token in the hypothesis was a button in the annotation interface and annotators were instructed to click on incorrect tokens to mark errors. Unmarked tokens were assumed OK. Additionally, they were instructed to keep markings minimal and only mark tokens that would be edited or deleted during post-editing.

¹For all Llama 13B experiments, we use a single Nvidia A40 GPU with 48GB VRAM on a shared server

²GPT-3.5-Turbo-0613 was used for all experiments involving OpenAI's GPT models in this paper

You will receive an annotation task called "Error Annotation"; the goal of this task is to mark each word in the machine translation as correct or incorrect. By default, all words are considered correct. By clicking on target words, they are marked as incorrect.

Here are the instructions in more detail:

You will be shown an English source sentence above and its machine translation into German below.

- Begin by reading the source sentence and then reading the translation.
- Consider which words would need to be deleted or changed in order to arrive at a correct translation.
- Mark the incorrect words of the translation by clicking on them.
- Clicking on the word causes a blue border to appear around the word. This word is now marked as incorrect.
- Clicking a second time will remove the blue border and it is now marked correct.
- Once all the incorrect words have a blue border, click on the "Next" button near the top of the page.
- Markings should be kept minimal. Mark only those terms that you would edit or delete in a post-editing scenario.
- If you would have to move a word to a different location, such as shifting a verb to the end of the sentence, mark it as incorrect.

If the translation contains no correct words or the source words are translated word by word but do not make sense together, mark them all as incorrect.

If the translation is correct as-is, proceed to the next annotation item.

If you cannot judge the quality of the translation because the source sentence is not comprehensible, or you are lacking domain knowledge to annotate wrong words, click the Skip button (to the right of the "Next" button) and then proceed to the next sentence.

The source sentences are taken from open-source software projects and documentation while the translations are produced by a generic machine translation system.

Figure 2: Instructions given to annotators on how to mark errors in sentences, including how to use the interface and desired marking behavior

Complete instructions are shown in Figure 2. Annotators could skip examples but must provide a reason. Reasons for skipping examples were "Source Incomprehensible", "Source Ambiguous", "Missing Knowledge", and "Other".

Annotation was split into two phases. Phase one was a trial run where all three annotators annotated the same 50 examples. In phase two, each annotator was given their own non-overlapping block of 500 source and hypothesis pairs. The phase one examples were used to compute summary statistics of annotation behavior, agreement coefficients, and to calibrate our instructions.

After phase two, filtering out skipped examples or those without any BAD markings yields a data set of 982 examples. We split this data set into two subsets; one set of size 492 for in-context examples and a set of 490 for test examples.

4.2 Annotation Statistics

Annotator 1 selected "Source Ambiguous" as the reason for skipping once and "Missing Knowledge" the other six times. Annotator 2 selected "Source Incomprehensible" for their skip. After removing the items skipped by any annotator, we have 43 examples that were annotated by all three.

Using the remaining common examples from phase one, we calculate the percentage of tokens marked per sentence and use that as a sentence-level quality judgment. This is then used to calculate Krippendorff's alpha (Krippendorff, 2004) to determine if our annotators agree on overall translation quality. We also calculate alpha on the token level OK/BAD annotations.

We calculated pair-wise Krippendorff's alpha in addition to the average agreement for both the sentence-level percentage marked and token-level annotations. The average amount of tokens marked for the unskipped sentences is visible in Table 1. Pairwise Krippendorff's alphas for percentage

Annotator	1	2	3
Percent Marked on Average	0.25	0.17	0.17
SD of Percent Marked	0.28	0.18	0.19

Table 1: Marking behaviors of each annotator in terms of percent of tokens marked in the trial annotation.

Annotator	2	3
1	0.258	0.481
2	\emptyset	0.222

Table 2: Inter-annotator agreement for percentage marked per sentence, given by Krippendorff’s Alpha.

Annotator	2	3
1	0.445	0.531
2	\emptyset	0.433

Table 3: Inter-annotator agreement for token classification, given by Krippendorff’s Alpha.

Annotator	1	2	3
Percent Marked on Average	0.10	0.19	0.09
SD of Percent Marked	0.10	0.15	0.1

Table 4: Marking behaviors of each annotator in terms of percent of tokens marked in the final annotation.

marked is visible in Table 2, while pairwise agreement for token classification is in Table 3. Average agreement for the percentage marked is $\alpha = 0.306$ and for token classification $\alpha = 0.466$. This suggests that, while agreement about overall sentence quality is not high, the reliability of classifying each token in the hypothesis is higher. These results were used to calibrate with the annotators after looking over the annotations made by each individual.

After calibration, we then assigned each annotator their block of 500 examples to annotate. Annotator 1 skipped 6 of the 500 sentences and annotator 2 skipped 20. Percentage marked was lower for annotators 1 and 3 during the full annotation as more sentences were left completely unmarked. Annotator 1 left 36% of sentences unmarked; annotator 2 left 23%; and annotator 3 left 38%. The percentage that was marked per sentence was also reduced after calibration, as shown in Table 4.

5 Experiments

5.1 Experimental Setup

Using the annotated data, we considered three machine translation tasks: Machine translation from scratch (*MT*); Automatic Post-editing (*APE*); and Post-Editing with error markings (*MRK*). Instructions were written for the LLM for each task and, for each example in the inference set, five examples were retrieved from the in-context example pool. We retrieve the most similar examples by using cosine similarity over SentenceTransformers (Reimers and Gurevych, 2019) embeddings computed on source sentences only³.

In *MT*, models were prompted to

Translate English to German.

and were shown five (source, reference) pairs. Full prompts can be found in the appendix A.1. For the *APE* task, models were prompted to

Read the English text and the German translation hypothesis and then correct the output. If the hypothesis is already correct, do not make any changes.

With this prompt, the models were given triples of (source, hypothesis, reference) with the hypothesis from our annotated data set and the reference coming from the parallel data.

In the *MRK* scenario, models were prompted to

Read the English text and the German translation hypothesis and then correct the output. Incorrect words are inside of tags '`<bad> </bad>`'. Please use this feedback in your correction. If the hypothesis is already correct, do not make any changes.

As with the *APE* prompt, models were given (source, hypothesis, reference) triples with the tokens that were marked as bad during annotation inside of XML-style tags, `<bad></bad>`. We decided that giving the error markings as in-line tags would be easier for the model to parse and integrate in its output than including another line where errors would be indicated further away from the corresponding tokens.

³We used the model *all-MiniLM-L6-v2* and retrieved the examples with the highest cosine similarity.

Condition	BLEU	TER	ME	UE	% Correct ME
Original Hyps	28.92	55.12	N.A.	N.A.	N.A.
MT (Llama/GPT)	29.83/38.61	55.97/49.21	N.A.	N.A.	N.A.
APE (Llama/GPT)	29.79/39.09	54.56/48.37	7.30/76.70	1.76/15.85	32% / N.A.
MRK (Llama/GPT)	30.09/39.31	54.70/48.32	14.76/78.36	3.60/13.90	67% / N.A.

Table 5: Results for both Llama 13B and GPT 3.5 across all metrics and translation scenarios (ME = Marking Edits, UE = Unmarking Edits, % Correct ME = Percentage of correct ME in manual evaluation).

5.2 Metrics

We evaluate the models’ new hypotheses with a suite of metrics to check for token level matches, semantic similarity, and error marking usage. We use the token based metrics BLEU⁴ (Papineni et al., 2002) and TER⁵ (Snover et al., 2006) as implemented in SacreBLEU (Post, 2018). We did not include the popular neural metric COMET (Rei et al., 2020) since is not sensitive to the individual token changes (Glushkova et al., 2023) that we ask the LLMs to perform.

In addition to these metrics, we also implement our own to see how well the models are at recognizing and making edits to errors. These are called *marking edit* (ME) and *unmarking edit* (UE). We perform a word-level diff in order to see which words need to be edited or deleted in the original hypothesis to arrive at the new hypothesis. Combining this with the error markings allows us to examine if edits were made to the tokens error-marked by the annotators (ME), or if tokens that were otherwise OK were changed (UE).

The ME and UE metrics, however, cannot tell if the edits were correct, only that they were made. To determine if the edits correctly fix the errors, we performed a manual evaluation on the marking edits produced by Llama 13B in both the *APE* and *MRK* settings. Three of the authors contributed to the evaluation. Two are native speakers of German and fluent in English while one is a native speaker of English and fluent in German. We selected 100 sentences with the most marking edits. Edits were evaluated in terms of their correctness, with a subjective yes or no answer given for the entire sentence.

⁴BLEU signature: nrefs:1 | case:mixed | eff:no | tok:13a | smooth:exp | version:2.4.0

⁵TER Signature nrefs:1 | case:lc | tok:tercom | norm:no | punct:yes | asian:no | version:2.4.0

6 Results

We show results across metrics for Llama 13B and GPT-3.5 in Table 5. Including error markings as input increases the frequency with which the models edits the marked tokens. For Llama 13B, we see editing rates for marked tokens double from 7.30 to 14.76. This suggests that, even after being asked to correct the hypotheses, Llama 13B finds its own outputs as acceptable translations. When errors are specifically pointed out to the model, it is much more capable of self-correcting errors.

Llama 13B nominally improves BLEU scores over the original hypotheses score (28.92) in all scenarios with *MRK* in the lead with 30.09, *MT* in second with 29.83⁶ and *APE* with 29.79. Nominal improvements over the original hypotheses are also found according to the TER metric, albeit only for *APE* and *MRK* scenarios.

The GPT model is already quite capable of finding errors in the hypotheses without error markings and the *APE* outputs achieve marking edits of 76.70 while *MRK* has a slight improvement of 78.36. Worth noting is the reduction in unmarking edits when prompting GPT with *MRK*. *MRK* reduces unmarking edits to 13.90 from 15.85 with *APE*. This means that indicating specific errors can constrain the number of edits that the GPT model makes. Additionally, nominal improvements of BLEU and TER scores are found in the *APE* and *MRK* scenarios over *MT* with GPT 3.5 as well. *MRK* improves BLEU to 39.31 from *MT*’s 38.61.

In the manual evaluation of marking edits, we found that *APE* made correct edits 32% of the time on average, while either making incorrect edits or not editing the rest. *MRK* on the other hand was judged correct 67% of the time on average. Agreement in terms of Krippendorff’s Alpha for sentence

⁶*MT* is able to surpass the original hypotheses with Llama 13B because the annotated hypotheses were generated with the same 5 examples for all inference segments while *MT* retrieved similar examples for each test segment.

level ratings of *APE* is $\alpha = 0.82$, while for *MRK* $\alpha = 0.55$. As *APE* makes fewer edits overall, it is easier to classify as incorrect or not editing. For *MRK* there was disagreement on how to handle partial edits or if not all markings were edited, requiring individual judgement by each evaluator.

7 Conclusion

We presented a pilot study to investigate the potential of augmenting a so-called PE-TM resource consisting of sources, machine translations, and human references, with human error markings in order to guide an LLM to self-correct marked erroneous term translations. We find that the LLM that produced the translation hypotheses identifies its own translations as correct, and therefore does not act on the instructions to correct errors. However, when prompted with error markings, the LLM learns to act on them, doubling the number of edits to marked tokens, with nearly 70% of the edits being correct according to a human evaluation. In sum, our pilot study shows that the additional effort of error marking a machine translation at test time allows an LLM translation system to learn focused corrections on marked errors from similar examples extracted from a PE-TM, leading to improved translation quality over *APE* and *MT*. In future work, we will investigate learned models for error markings. These require larger TMs for reliable training of markings estimators, but also bear the promise of improved retrieval augmentation.

8 Acknowledgements

The second author acknowledges support by the state of Baden-Württemberg through bwHPC and the German Research Foundation (DFG) through grant INST 35/1597-1 FUGG.

References

- Berger, Nathaniel, Miriam Exel, Matthias Huck, and Stefan Riezler. 2023. Enhancing supervised learning with contrastive markings in neural machine translation training. In *Proceedings of the 24th Annual Conference of The European Association for Machine Translation (EAMT)*, Tampere, Finland.
- Domingo, Miguel, Álvaro Peris, and Francisco Casacuberta. 2017. Segment-based interactive-predictive machine translation. *Machine Translation*, 31(4):163–185.
- Exel, Miriam, Bianka Buschbeck, Lauritz Brandt, and Simona Doneva. 2020. Terminology-constrained neural machine translation at SAP. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation (EAMT)*, Lisboa, Portugal.
- Garcia, Xavier, Yamini Bansal, Colin Cherry, George Foster, Maxim Krikun, Melvin Johnson, and Orhan Firat. 2023. The unreasonable effectiveness of few-shot learning for machine translation. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, Hanolulu, Hawaii, USA.
- Ghazvininejad, Marjan, Hila Gonen, and Luke Zettlemoyer. 2023. Dictionary-based phrase-level prompting of large language models for machine translation. *arXiv*, abs/2302.07856.
- Glushkova, Taisiya, Chrysoula Zerva, and André F. T. Martins. 2023. BLEU meets COMET: Combining lexical and neural metrics towards robust machine translation evaluation. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation (EACL)*, Tampere, Finland.
- Grangier, David and Michael Auli. 2018. QuickEdit: Editing text & translations by crossing words out. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL:HLT)*, New Orleans, Louisiana.
- Hoang, Cuong, Devendra Sachan, Prashant Mathur, Brian Thompson, and Marcello Federico. 2023. Improving retrieval augmented neural machine translation by controlling source and fuzzy-match interactions. In *Findings of the Association for Computational Linguistics: EACL 2023*, Dubrovnik, Croatia.
- Huang, Jie, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2023. Large language models cannot self-correct reasoning yet. *arXiv*, abs/2310.01798.
- Joulin, Armand, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, Valencia, Spain.
- Kreutzer, Julia, Nathaniel Berger, and Stefan Riezler. 2020. Correct me if you can: Learning from error corrections and markings. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation (EAMT)*, Lisbon, Portugal.
- Krippendorff, Klaus. 2004. Reliability in content analysis: Some common misconceptions and recommendations. *Human Communication Research*, 30(3):411–433.
- Lam, Tsz Kin, Shigehiko Schamoni, and Stefan Riezler. 2019. Interactive-predictive neural machine translation through reinforcement and imitation. In *Proceedings of the Machine Translation Summit (MT-SUMMIT XVII)*, Dublin, Ireland.

- Madaan, Aman, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Sean Welleck, Bodhisattwa Prasad Majumder, Shashank Gupta, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. *arXiv*, abs/2303.17651.
- Marie, Benjamin and Aurélien Max. 2015. Touch-based pre-post-editing of machine translation output. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Lisbon, Portugal.
- Moslem, Yasmin, Rejwanul Haque, John D. Kelleher, and Andy Way. 2023. Adaptive machine translation with large language models. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation (EAMT)*, Tampere, Finland.
- OpenAI. 2023. GPT-4 technical report. Technical report, OpenAI.
- Pan, Liangming, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 2023. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. *arXiv*, abs/2308.03188.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics (ACL)*, Philadelphia, PA.
- Petrushkov, Pavel, Shahram Khadivi, and Evgeny Matusov. 2018. Learning from chunk-based feedback in neural machine translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)*, Melbourne, Australia.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the 3rd Conference on Machine Translation (WMT)*, Brussels, Belgium.
- Raunak, Vikas, Amr Sharaf, Yiren Wang, Hany Hassan Awadallah, and Arul Menezes. 2023. Leveraging GPT-4 for automatic translation post-editing. In *Findings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Singapore.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Online.
- Reimers, Nils and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Hong Kong, China.
- Snover, Matthew, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas (AMTA'06)*, Cambridge, MA.
- Tiedemann, Jörg. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, Istanbul, Turkey.
- Touvron, Hugo, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and efficient foundation language models. *arXiv*, abs/2302.13971.
- Turchi, Marco, Matteo Negri, M. Amin Farajian, and Marcello Federico. 2017. Continuous learning from human post-edits for neural machine translation. *The Prague Bulletin of Mathematical Linguistics (PBML)*, 1(108):233–244.
- Vidal, Blanca, Albert Llorens, and Juan Alonso. 2022. Automatic post-editing of MT output using large language models. In Campbell, Janice, Stephen Larocca, Jay Marciano, Konstantin Savenkov, and Alex Yanishevsky, editors, *Proceedings of the 15th Biennial Conference of the Association for Machine Translation in the Americas (AMTA)*, Orlando, Florida, USA.
- Vilar, David, Markus Freitag, Colin Cherry, Jiaming Luo, Viresh Ratnakar, and George Foster. 2023. Prompting PaLM for translation: Assessing strategies and performance. In Rogers, Anna, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL)*, Toronto, Canada.
- Wolf, Thomas, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Online.
- Xu, Wenda, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, Biao Zhang, Zhongtao Liu, William Yang Wang, Lei Li, and Markus Freitag. 2023. Pinpoint, not criticize: Refining large language models via fine-grained actionable feedback.

A Appendix

A.1 Example Prompts

An example of a prompt for *MT*, *APE*, and *MRK* are in Figures 3, 4, and 5, respectively.

Translate English to German.

English: Cookies are part of the HTTP header, so `setcookie()` must be called before any output is sent to the browser.

German: Sie sind Bestandteil des HTTP-Headers, was bedeutet, dass die Funktion `setcookie()` aufgerufen werden muss, bevor irgendeine Ausgabe an den Browser erfolgt.

English: `session.use_only_cookies` specifies whether the module will only use cookies to store the session id on the client side.

German: `session.use_only_cookies` spezifiziert, ob das Modul nur Cookies verwendet, um die Session-ID clientseitig zu speichern.

English: Note that `SID` is only defined if the client didn't send the right cookie.

German: Beachten Sie, dass `SID` nur definiert ist, wenn vom Client nicht das richtige Cookie gesendet wurde.

English: The server does not support the request type of the body.

German: Der Server unterstützt den angeforderten Typ nicht.%1: request type

English: Must be in active session on local console

German: Nur in aktiver Sitzung auf lokaler Konsole

English: Like other headers, cookies must be sent before any output from your script (this is a protocol restriction).

German:

Figure 3: Example of 5-shot prompt for English-to-German Translation. Each demonstration example consists of a source segment in English (in green), and a reference translation (in red).

Read the English text and the German translation hypothesis and then correct the output. If the hypothesis is already correct, do not make any changes.

English: Cookies are part of the HTTP header, so `setcookie()` must be called before any output is sent to the browser.

Hypothesis: Cookies sind Teil des HTTP-Headers , deshalb muss `setcookie()` vor jedem Ausgabe-Output an den Browser aufgerufen werden .

German: Sie sind Bestandteil des HTTP-Headers, was bedeutet, dass die Funktion `setcookie()` aufgerufen werden muss, bevor irgendeine Ausgabe an den Browser erfolgt.

English: `session.use_only_cookies` specifies whether the module will only use cookies to store the session id on the client side.

Hypothesis: `session.use_only_cookies` bestimmt , ob das Modul nur mit Cookies die Session-ID auf dem Client-Betriebssystem speichert .

German: `session.use_only_cookies` spezifiziert, ob das Modul nur Cookies verwendet, um die Session-ID clientseitig zu speichern.

English: Note that SID is only defined if the client didn't send the right cookie.

Hypothesis: Beachtet , dass SID nur definiert ist , wenn der Client nicht den richtigen Cookie gesendet hat .

German: Beachten Sie, dass SID nur definiert ist, wenn vom Client nicht das richtige Cookie gesendet wurde.

English: The server does not support the request type of the body.

Hypothesis: Der Server unterstützt nicht die Anforderungstyp der Body .

German: Der Server unterstützt den angeforderten Typ nicht.%1: request type

English: Must be in & active session on local console

Hypothesis: Muss in & aktiver Sitzung auf dem lokalen Konsole

German: Nur in & aktiver Sitzung auf lokaler Konsole

English: Like other headers, cookies must be sent before any output from your script (this is a protocol restriction).

Hypothesis: Wie andere Headern müssen Cookies vor jedem Ausgabe-Output (dies ist eine Protokoll-Einschränkung) gesendet werden .

German:

Figure 4: Example of 5-shot prompt for English-to-German Automatic Post-Editing (APE). Each demonstration example consists of a source segment in English (in green), a translation hypothesis in German (in blue), and a reference translation (in red).

Read the English text and the German translation hypothesis and then correct the output. Incorrect words are inside of tags '**<bad>** **</bad>**'. Please use this feedback in your correction. If the hypothesis is already correct, do not make any changes.

English: Cookies are part of the HTTP header, so `setcookie()` must be called before any output is sent to the browser.

Hypothesis: Cookies sind Teil des HTTP-Headers , deshalb muss `setcookie()` vor jedem **<bad>** Ausgabe-Output **</bad>** an den Browser aufgerufen werden .

German: Sie sind Bestandteil des HTTP-Headers, was bedeutet, dass die Funktion `setcookie()` aufgerufen werden muss, bevor irgendeine Ausgabe an den Browser erfolgt.

English: `session.use_only_cookies` specifies whether the module will only use cookies to store the session id on the client side.

Hypothesis: `session .use_only_cookies` bestimmt , ob das Modul nur **<bad>** mit **</bad>** Cookies die Session-ID auf dem **<bad>** Client-Betriebssystem speichert **</bad>** .

German: `session.use_only_cookies` spezifiziert, ob das Modul nur Cookies verwendet, um die Session-ID clientseitig zu speichern.

English: Note that SID is only defined if the client didn't send the right cookie.

Hypothesis: **<bad>** Beachtet **</bad>** , dass SID nur definiert **<bad>** ist **</bad>** , wenn der Client nicht den richtigen Cookie gesendet hat .

German: Beachten Sie, dass SID nur definiert ist, wenn vom Client nicht das richtige Cookie gesendet wurde.

English: The server does not support the request type of the body.

Hypothesis: Der Server unterstützt nicht **<bad>** die **</bad>** Anforderungstyp **<bad>** der **</bad>** Body .

German: Der Server unterstützt den angeforderten Typ nicht.%1: request type

English: Must be in & active session on local console

Hypothesis: Muss in & aktiver Sitzung auf **<bad>** dem **</bad>** lokalen Konsole

German: Nur in & aktiver Sitzung auf lokaler Konsole

English: Like other headers, cookies must be sent before any output from your script (this is a protocol restriction).

Hypothesis: Wie andere **<bad>** Headern **</bad>** müssen Cookies vor jedem **<bad>** Ausgabe-Output **</bad>** (dies ist eine Protokoll-Einschränkung) gesendet werden .

German:

Figure 5: Example of 5-shot prompt for English-to-German Post-Editing with error markings (*MRK*). Error markings inside by tags **<bad>** **</bad>**. Each demonstration example consists of a source segment in English (in green), a translation hypothesis in German (in blue), and a reference translation (in red).

Estonian-Centric Machine Translation: Data, Models, and Challenges

Elizaveta Korotkova and Mark Fishel

Institute of Computer Science

University of Tartu, Estonia

{elizaveta.korotkova, mark.fisel}@ut.ee

Abstract

Machine translation (MT) research is most typically English-centric. In recent years, massively multilingual translation systems have also been increasingly popular. However, efforts purposefully focused on less-resourced languages are less widespread. In this paper, we focus on MT from and into the Estonian language. First, emphasizing the importance of data availability, we generate and publicly release a back-translation corpus of over 2 billion sentence pairs. Second, using these novel data, we create MT models covering 18 translation directions, all either from or into Estonian. We re-use the encoder of the NLLB multilingual model and train modular decoders separately for each language, surpassing the original NLLB quality. Our resulting MT models largely outperform other open-source MT systems, including previous Estonian-focused efforts, and are released as part of this submission.

1 Introduction

The majority of work on neural machine translation (NMT) is nowadays primarily English-centric, with some notable work on (massively) multilingual MT (Fan et al., 2020; NLLB Team et al., 2022; Kudugunta et al., 2023). In recent years, some attention has been directed at translation directions out of English (e.g. this is the primary focus of the WMT’2024 evaluation campaign¹) or at

© 2024 The authors. This article is licensed under a Creative Commons 4.0 licence, no derivative works, attribution, CC-BY-ND.

¹<http://www2.statmt.org/wmt24/translation-task.html>

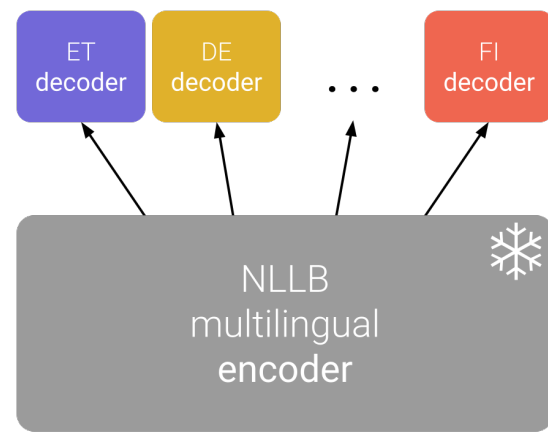


Figure 1: Model architecture. We reuse the multilingual Transformer encoder of NLLB-1.3B and train a new Transformer decoder for each target language.

pairs that do not include English: for instance, recent WMT and IWSLT shared tasks included one or two such pairs (Kocmi et al., 2023; Kocmi et al., 2022; Agarwal et al., 2023).

In this work, we present our recent efforts on advancing Estonian-centric machine translation. In a broader scope the work is part of the Neurotõlge project, which develops open machine translation for Estonian.² The name Neurotõlge means *Neural translation* in Estonian and the work on its development has started in 2017 and is ongoing.

The present contribution covers 18 new translation directions for Neurotõlge from and into Estonian. We openly release a massive back-translation corpus for these language pairs, extending the Synthetic Corpus of Parallel Estonian (SynEst) (Korotkova et al., in press), and release translation models trained using these data.

We employ a partially modular approach (Escolano et al., 2021; Lyu et al., 2020) in creating

²<https://translate.ut.ee>

translation models. Specifically, we use the encoder of an existing massively multilingual translation system NLLB (NLLB Team et al., 2022) and create the decoders for each target language as separate modules (the architecture is shown in Figure 1). This setup makes it possible to train the decoders independently, and any subset of the decoders can be deployed afterwards. The achieved translation quality is better than the original NLLB system and also surpasses other open systems on the included translation directions.

The main contributions of this paper are thus:

- we extend the SynEst corpus to cover 12 new translation directions and 4 new data sources, adding over 2 billion filtered sentence pairs to the corpus, and make the full corpus publicly available;³
- we create new MT systems for Estonian translation, covering 6 translation directions from Estonian and 12 translation directions into Estonian. Our systems demonstrate stronger translation performance than previous open-source efforts, including Estonian-centric ones, on most language pairs when translating from Estonian into other languages, and show especially noticeable and consistent improvements for translation into Estonian (up to 13 BLEU (Papineni et al., 2002) depending on translation direction and text domain). The models are released for open use.⁴

2 Related Work

In our work, we focus on strengthening the capabilities of open-source MT systems focused on the Estonian language. This builds upon previous efforts centered on Estonian public translation, most recently, the MTEE governmental project (Tättar et al., 2022), and, more generally, the Neurotõlge project and online translation engine.² MTEE covered translation between Estonian and three other languages: English, German, and Russian, and achieved state-of-the-art translation quality at the time (Tättar et al., 2022). In this work, we train

³<https://metashare.ut.ee/repository/search/?q=SynEst>, for direct DOI links to each language pair, see Appendix B.

⁴<https://huggingface.co/tartuNLP/synest-models>

Estonian-centric models for more language pairs, outperforming the MTEE models in most cases.

Instead of training models from scratch, we use the NLLB multilingual translation model (NLLB Team et al., 2022) as a starting point for our systems. NLLB is a massive effort utilizing the multilingual MT approach (Dong et al., 2015; Johnson et al., 2017), and covering 200 languages, which makes it a convenient base on which to build systems tailored to a smaller number of languages.

In this work, we mostly rely on creating large amounts of new training data to improve Estonian translation. Specifically, we use the back-translation technique (Sennrich et al., 2016). Existing MT systems are used to generate translations of monolingual corpora into desired languages. The obtained parallel data is then reversed and used to augment the training corpus. Thus, the noisy, automatically translated text is on the source side, and the target side contains the cleaner original data, which allows the model to learn text generation based on genuine data. Specifically, we use and extend the SynEst corpus (Korotkova et al., in press), an Estonian-focused back-translation dataset, to cover new translation directions and source corpora.

In terms of model architecture, our systems are inspired by modular approaches (Lyu et al., 2020; Escolano et al., 2021), where multilingual MT models share encoder and decoder modules for each input and output language instead of having one encoder and one decoder covering all languages. More specifically, we use an existing multilingual encoder module from NLLB and train a new decoder for each target language from scratch, somewhat similarly to concurrent work on "mix-and-match translation" by Purason et al. (2024), where encoders and decoders from different models are unified to form a new model.

3 Extending the SynEst Corpus

Synthetic Corpus of Parallel Estonian, or SynEst (Korotkova et al., in press), includes data from the NewsCrawl monolingual corpus (Kocmi et al., 2023) automatically translated into Estonian from 11 languages (Arabic, Chinese, English, Finnish, French, German, Latvian, Lithuanian, Russian, Spanish, and Ukrainian). The dataset can be used as a back-translated corpus to facilitate training MT models which include Estonian.

In this work, we significantly extend SynEst to

code	target language	parallel	back-translated corpus				total
			NewsCrawl	ParaCrawl	UNPC	OpenSubtitles	
DE	German	9.3	332.6	159.3	–	–	501.2
EN	English	19.6	254.7	433.2	19.4	61.0	787.9
FI	Finnish	15.0	23.5	19.1	–	–	57.6
RU	Russian	5.1	86.8	2.2	13.1	–	107.2
UK	Ukrainian	2.6	1.8	6.7	–	–	11.1
ZH	Chinese	5.8	10.4	4.7	–	–	20.9

Table 1: Sizes of training corpora for models translating from Estonian into other languages (filtered, in millions of sentence pairs). Parallel shows the total size of all parallel corpora used for each language pair. For back-translated corpora, the source side (Estonian) is the automatically translated data, while the target side is the original data. UNPC denotes the United Nations Parallel Corpus.

code	source language	parallel	ENC	total
AR	Arabic	6.3	94.3	100.6
DE	German	9.3	143.8	153.1
EN	English	19.6	144.7	164.3
ES	Spanish	19.5	126.8	146.3
FI	Finnish	15.0	136.8	151.8
FR	French	18.8	132.1	150.9
LT	Lithuanian	10.5	132.7	143.2
LV	Latvian	7.1	132.2	139.3
RU	Russian	5.1	112.1	117.2
SV	Swedish	13.4	127.8	141.2
UK	Ukrainian	2.6	115.7	118.3
ZH	Chinese	5.8	113.6	119.4
total				1,645.6

Table 2: Sizes of training corpora for models translating into Estonian from other languages (filtered, in millions of sentence pairs). Parallel shows the total size of all parallel corpora used for each language pair. ENC denotes the Estonian Parallel Corpus. The Estonian Parallel Corpus was back-translated: the source side is the data automatically translated from Estonian into other languages, while the target side is the original Estonian data.

include more source corpora and translation directions, most importantly, introducing translation directions *from* Estonian. We make the updated dataset publicly available for unrestricted use.³

3.1 Translation Directions into Estonian

For translation directions into Estonian, we extend the corpus with three new data sources: ParaCrawl (Bañón et al., 2020), the United Nations Parallel Corpus (Ziemski et al., 2016), and OpenSubtitles (Lison and Tiedemann, 2016).

In case of ParaCrawl, we use 10 language pairs present in this parallel corpus: one side is al-

ways English, and the other one of German, Spanish, Finnish, French, Lithuanian, Latvian, Russian, Swedish, Ukrainian, and Chinese. We automatically translate both sides of the corpora into Estonian. The sizes of the resulting corpora range from 5.4 million sentence pairs for Russian–Estonian to a total of 878.4 million pairs for English–Estonian. As both sides of the parallel corpus are translated into a third language (Estonian), this setup opens the possibility of exploring triangular MT approaches; however at present we treat the corpora we translate as monolingual and leave investigation of this direction for future work.

For the United Nations Parallel Corpus, we translate its English and Russian monolingual subsets into Estonian, obtaining 33.4 million and 28.5 million sentence pairs before filtering, respectively. Finally, we translate the English OpenSubtitles corpus into Estonian as well, resulting in 441.4 million sentence pairs before filtering.

The total sizes of the generated dataset for each source corpus and translation direction are given in Table 8 in Appendix A.

3.2 Translation Directions from Estonian

Most importantly, we focus on extending the SynEst synthetic corpus to include translation directions from Estonian. This will allow to use the corpus to train models for translation into Estonian. We translate the Estonian National Corpus (Koppel and Kallas, 2022) into 12 languages: Arabic, Chinese, English, Finnish, French, German, Latvian, Lithuanian, Russian, Spanish, Swedish, and Ukrainian. The resulting back-translation corpus contains between 171.4 million and 196.6 million sentence pairs per translation direction (see Table 7 in Appendix A for approximate numbers

	target language					
	DE	EN	FI	RU	UK	ZH
NLLB-1.3B	24.4	36.7	15.5	22.4	18.7	25.0
MTEE	25.8	37.0	–	22.4	–	–
MADLAD-3B	26.0	37.8	20.1	20.0	15.5	<u>33.5</u>
Ours	<u>27.5</u>	<u>38.1</u>	<u>21.9</u>	<u>23.5</u>	<u>21.3</u>	31.6
DeepL	30.9	39.9	24.4	26.7	25.6	40.5
Google	30.8	41.7	22.9	26.6	24.4	42.2

Table 3: BLEU scores on the FLORES-devtest benchmark for models translating from Estonian into other languages. The best scores overall are shown in **bold**, and the best scores among open-source models are underlined. For our models, we report the score of the checkpoint with the best validation BLEU. With MTEE, we use the general-domain model to translate the FLORES benchmark.

for each translation direction).

3.3 Translation Models

For generating the synthetic side of the SynEst corpus we translate from and into English, German, and Russian with the MTEE models (Tättar et al., 2022), using the domain-specific engines MTEE-legal for the United Nations Parallel Corpus and MTEE-general for all other corpora. For translation directions not involving these languages we use the M2M-100 1.2B-parameter model (Fan et al., 2020). In all cases, we use beam search with beam size 5.

4 Experiments

4.1 Models

We replicate the model setup used in previous exploratory experiments (Korotkova et al., in press). We base our systems on the multilingual NLLB-1.3B dense model (NLLB Team et al., 2022). We freeze the NLLB encoder and train a new, randomly initialized Transformer decoder (Vaswani et al., 2017) for each target language. We keep the dimensions of the decoder layers the same as in the encoder, but use 6 decoder layers instead of the encoder’s 24. Keeping the encoder parameters fixed allows to reduce the training-time costs, while reducing the size of the decoder lowers both training- and inference-time costs compared to full fine-tuning of the base model. Freezing the encoder parameters also maintains the multilingual properties of the encoder, meaning that after fine-tuning the model on a certain translation direction it can still translate from any of the 200 languages of NLLB. As all models share the same encoder parameters, final models can be built in a modular

fashion, with a single decoder for all translation directions, and one encoder per target language.

We focus on creating Estonian-centric MT models: all translation directions in our experiments include Estonian as either the source or the target language. Specifically, for translation from Estonian into other languages, we train models that translate into German, English, Finnish, Russian, Ukrainian, and Chinese. For translation into Estonian, as the encoder is shared between all models and Estonian is the common target language, we train a single model on the concatenation of data representing 12 language pairs (see Table 2).

We use FairSeq (Ott et al., 2019) to train our models; details on model and training hyperparameters can be found in Appendix D.

4.2 Training Data

To train our models, we use two types of data: parallel corpora and the extended SynEst back-translated corpus.

We use the concatenation of 10 parallel corpora: CCMatrix (Schwenk et al., 2021b), WikiMatrix (Schwenk et al., 2021a), MultiParaCrawl (Bañón et al., 2020), Europarl (Koehn, 2005), OpenSubtitles (Lison and Tiedemann, 2016), JRC-Acquis (Steinberger et al., 2006), TED2020 (Reimers and Gurevych, 2020), EMEA, infopankki, and DGT (Tiedemann, 2012). For the Estonian–English language pair, MultiParaCrawl is replaced with ParaCrawl (Bañón et al., 2020). Not all of these corpora exist for each language pair in our experiments; we use each of the corpora whenever it is available for a language pair.

For SynEst, we use all source corpora available for a given translation direction. As the dataset is used as additional back-translation data, the auto-

	ET-DE	ET-EN	ET-RU
News			
NLLB-1.3B	25.8	25.6	22.8
MTEE	30.1	26.4	26.9
MADLAD-3B	26.3	<u>28.7</u>	19.7
Ours	30.5	25.9	26.5
DeepL	28.0	28.1	23.5
Google	26.0	30.0	21.2
Crisis			
NLLB-1.3B	26.3	21.4	26.2
MTEE	29.8	33.8	33.8
MADLAD-3B	22.1	<u>35.0</u>	25.0
Ours	30.3	33.2	34.7
DeepL	28.1	34.1	27.3
Google	26.6	36.1	27.6
Military			
NLLB-1.3B	21.0	31.1	30.1
MTEE	24.2	35.4	35.9
MADLAD-3B	19.6	33.2	28.8
Ours	25.4	32.9	35.7
DeepL	20.0	32.7	31.0
Google	20.3	34.2	34.5
Legal			
NLLB-1.3B	27.1	48.9	35.5
MTEE	34.0	55.1	42.8
MADLAD-3B	32.1	47.8	39.9
Ours	<u>34.7</u>	53.7	43.0
DeepL	34.8	50.9	35.5
Google	39.1	50.9	37.8
Spoken			
NLLB-1.3B	29.3	30.5	23.3
MTEE	33.0	34.3	28.1
MADLAD-3B	33.1	<u>35.2</u>	22.8
Ours	<u>33.2</u>	32.2	28.0
DeepL	29.9	34.4	23.5
Google	36.0	41.0	22.3

Table 4: BLEU scores on the MTEE domain benchmark sets for models translating from Estonian into other languages. The best scores overall are shown in **bold**, and the best scores among open-source models are underlined. For our models, we report the score of the checkpoint with the best validation BLEU. With MTEE, we show the scores reported by Tättar et al. (2022).

matically generated side of the corpus is always used as the source and the cleaner original data as

the target during training.

We concatenate all corpora to create our full training dataset. Approximate sizes of the full training corpora and their components are shown in Tables 1 and 2 for model translation directions from Estonian and into Estonian, respectively. (The sizes are shown after filtering; details on data filtering can be found in Appendix C).

The dev split of the FLORES dataset (Goyal et al., 2022) is used as the validation set.

4.3 Evaluation

We compare the performance of our Estonian-centric models to that of three other open-source MT systems:

- the NLLB-1.3B (NLLB Team et al., 2022) multilingual translation model, which also serves as the starting model in our experiments;
- the models trained within the MTEE project (Tättar et al., 2022), which was the previous effort of public Estonian-centric MT. These models cover the Estonian↔German, Estonian↔English, and Estonian↔Russian translation directions, and employ a fully modular approach;
- the more recent MADLAD-400 3B (Kudugunta et al., 2023).

For additional comparison, we also show the results of DeepL⁵ and Google Translate,⁶ two widely used proprietary online translation engines.

The test sets we employ for evaluation are the FLORES evaluation benchmark (Goyal et al., 2022) (the devtest split), and the MTEE domain-specific benchmark sets (Tättar et al., 2022). FLORES is useful in providing a benchmark for multilingual translation between many languages, which is based on Wikipedia. MTEE, while covering fewer language pairs (Estonian–English, Estonian–German, and Estonian–Russian), is centered on language pairs which include Estonian, and allows to estimate model performance on text belonging to 5 distinct domains.

We use the sacreBLEU implementation (Post, 2018) of the BLEU score (Papineni et al., 2002) to

⁵<https://www.deepl.com/translator>

⁶<https://translate.google.com>

	source language											
	AR	DE	EN	ES	FI	FR	LT	LV	RU	SV	UK	ZH
NLLB-1.3B	15.7	17.8	22.7	13.8	16.1	17.3	15.1	16.1	15.8	18.4	16.9	11.6
MTEE	–	21.7	27.6	–	–	–	–	–	<u>20.2</u>	–	–	–
MADLAD-3B	<u>20.3</u>	21.7	26.2	16.3	19.2	19.9	<u>19.3</u>	<u>22.8</u>	17.7	21.3	16.2	<u>15.4</u>
Ours	20.0	<u>23.0</u>	<u>29.4</u>	<u>16.7</u>	<u>20.9</u>	<u>23.3</u>	<u>19.3</u>	21.0	20.1	<u>24.0</u>	<u>21.4</u>	14.6
DeepL	23.4	24.4	30.2	19.0	22.5	23.7	22.1	23.6	22.6	26.3	24.1	18.0
Google	23.2	25.3	30.7	18.5	22.4	24.5	21.5	23.6	22.6	25.7	23.3	18.8

Table 5: BLEU scores on the FLORES-devtest benchmark for models translating from other languages into Estonian. The best scores overall are shown in **bold**, and the best scores among open-source models are underlined. For our model, we report the score of the checkpoint with the best validation loss (the same checkpoint is used for all source languages). With MTEE, we use the general-domain model to translate the FLORES benchmark.

measure the models’ performance.⁷ Additionally, we report COMET scores (Rei et al., 2020) in Appendix E. For models translating from Estonian, we choose the checkpoint which shows the best BLEU score on FLORES-dev for the language pair in question. For the models translating into Estonian, we use the checkpoint showing the best loss on the combined validation set; we do not choose a best checkpoint for each source language separately.

5 Results

BLEU scores of NLLB-1.3B, MTEE, MADLAD-3B, our model, DeepL, and Google Translate on FLORES-devtest for translation directions from Estonian into other languages (our experiments cover German, English, Finnish, Russian, Ukrainian, and Chinese as target languages) are shown in Table 3. In this setting, our model shows the strongest results among the open-source systems for five out of six language pairs, outperforming the next best open-source models by 0.3 to 2.6 BLEU points. On the MTEE domain benchmarks (Table 4), our model consistently outperforms other open-source ones on the Estonian–German language pair, while for Estonian–English it shows lower scores than the MTEE and, for most domains, MADLAD models. For Estonian–Russian, results are more mixed, with our models being the best among all models on the crisis and legal domains (with a small margin of 0.2 BLEU over MTEE for legal and a more noticeable one of 0.9 BLEU for crisis) and falling slightly behind MTEE on the news,

⁷sacreBLEU signature for all target languages except Chinese: nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1. For Chinese: the same with tok:zh.

military, and spoken domains (by up to 0.4 BLEU).

Table 5 shows results on FLORES-devtest for translation into Estonian. Our model noticeably improves upon the NLLB model for all translation directions, while also outperforming all compared open-source models on 7 out of 12 translation directions. On three more directions, the difference between our model and the best performing one among open systems does not exceed 0.3 BLEU points.

From Table 6 we see that our into-Estonian model performs consistently well on different domains. It outperforms all models, including proprietary ones and the MTEE models fine-tuned to these domains, on all language pairs and domains, with the exception of EN–ET news, with margins to the next best models ranging from 0.2 to 13 BLEU for different language pairs and domains. This consistently strong performance can be attributed to the fact that this single model has encountered a vast amount of training data, with 12 input languages and Estonian as the output language, leading it to learn generating Estonian output very well.

6 Deployment and Known Issues

The models are made publicly available on the HuggingFace model hub⁴ and can be run using the TartuNLP translation worker.⁸ The models are set up in a modular fashion, with one encoder covering all input languages and a separate decoder for each output language.

We have found that the models are not robust to some inputs, such as single words; while full

⁸<https://github.com/TartuNLP/translation-worker/tree/nllb-based-est>

	DE-ET	EN-ET	RU-ET
News			
NLLB-1.3B	22.0	15.6	19.5
MTEE	29.7	18.0	27.2
MADLAD-3B	24.9	19.0	22.5
Ours	<u>33.2</u>	<u>19.7</u>	<u>30.0</u>
DeepL	29.5	21.4	23.0
Google	28.9	19.7	24.8
Crisis			
NLLB-1.3B	27.4	24.3	20.1
MTEE	40.1	41.6	38.4
MADLAD-3B	36.2	31.1	27.2
Ours	<u>53.1</u>	<u>45.8</u>	<u>40.8</u>
DeepL	38.7	37.2	28.8
Google	39.6	41.2	32.3
Military			
NLLB-1.3B	22.6	21.6	20.1
MTEE	31.9	30.2	30.8
MADLAD-3B	28.0	24.6	24.1
Ours	<u>37.1</u>	<u>31.9</u>	<u>32.7</u>
DeepL	31.2	31.7	26.2
Google	28.6	31.7	26.8
Legal			
NLLB-1.3B	25.0	31.1	26.9
MTEE	32.4	50.8	47.1
MADLAD-3B	31.1	31.7	37.9
Ours	<u>48.0</u>	<u>52.1</u>	<u>50.3</u>
DeepL	39.2	47.8	37.0
Google	37.4	48.7	38.7
Spoken			
NLLB-1.3B	23.0	18.0	16.9
MTEE	31.7	23.7	24.4
MADLAD-3B	27.5	22.2	19.5
Ours	<u>37.5</u>	<u>26.1</u>	<u>27.3</u>
DeepL	30.7	24.2	19.1
Google	27.9	23.6	19.2

Table 6: BLEU scores on the MTEE domain benchmark sets for models translating from other languages into Estonian. The best scores overall are shown in **bold**, and the best scores among open-source models are underlined. For our model, we report the score of the checkpoint with the best validation loss (the same checkpoint is used for all source languages). With MTEE, we show the scores reported by Tättar et al. (2022).

sentence translation works reasonably well, with single-word or isolated phrase input the models

may start severely overgenerating.

7 Future Work

So far the efforts of the project have focused on sentence-level NMT. The next iterations of development and model training will likely focus on document-level MT, either with sequence-to-sequence or decoder-only models. Moreover, we are looking into instruction-tuned sequence-to-sequence models: this approach should yield translation-specific emergent abilities and would thus enable the integration of terminologies, on-the-fly domain adaptation, and other types of translation output control. We also plan to dedicate more attention to the robustness of the developed translation engines, for instance, by including upper-cased data in the training dataset for smoother handling of headlines and other all-caps segments, as well as including phrase and word pairs to enhance translation performance when the input is not a complete sentence.

8 Conclusion

In this work, we have made a contribution towards open-source machine translation centered on the Estonian language.

First, we presented an extended version of the SynEst synthetic corpus. The new version introduces 12 translation directions from Estonian, in addition to previously present directions into Estonian. In total, we have generated over 2 billion filtered sentence pairs. We release the full corpus for public use and hope that the availability of this resource will facilitate further work on Estonian translation.

Second, we created new MT models for translation from Estonian into 6 languages and from 12 languages into Estonian and made them publicly available. Evaluation on two benchmarks covering 6 domains has shown that our models are comparable to or outperform previous open efforts on translation from Estonian, depending on the language pair and domain, and perform especially well on translation into Estonian, outperforming not only previous open-source but also proprietary systems by up to 13 BLEU on some domains. These consistent improvements are likely due to the use of massive amounts of synthetic data we created.

References

- Agarwal, Milind, Sweta Agrawal, Antonios Anastasopoulos, Luisa Bentivogli, Ondřej Bojar, Claudia Borg, Marine Carpuat, Roldano Cattoni, Mauro Cettolo, Mingda Chen, William Chen, Khalid Choukri, Alexandra Chronopoulou, Anna Currey, Thierry Declerck, Qianqian Dong, Kevin Duh, Yannick Estève, Marcello Federico, Souhir Gahbiche, Barry Haddow, Benjamin Hsu, Phu Mon Htut, Hirofumi Inaguma, Dávid Javorský, John Judge, Yasumasa Kano, Tom Ko, Rishu Kumar, Pengwei Li, Xutai Ma, Prashant Mathur, Evgeny Matusov, Paul McNamee, John P. McCrae, Kenton Murray, Maria Nadejde, Satoshi Nakamura, Matteo Negri, Ha Nguyen, Jan Niehues, Xing Niu, Atul Kr. Ojha, John E. Ortega, Proyag Pal, Juan Pino, Lonneke van der Plas, Peter Polák, Elijah Rippeth, Elizabeth Salesky, Jiatong Shi, Matthias Sperber, Sebastian Stüker, Katsuhito Sudoh, Yun Tang, Brian Thompson, Kevin Tran, Marco Turchi, Alex Waibel, Mingxuan Wang, Shinji Watanabe, and Rodolfo Zevallos. 2023. Findings of the IWSLT 2023 evaluation campaign. In *Proceedings of the 20th International Conference on Spoken Language Translation (IWSLT 2023)*, pages 1–61, Toronto, Canada (in-person and online).
- Aulamo, Mikko, Sami Virpioja, and Jörg Tiedemann. 2020. OpusFilter: A configurable parallel corpus filtering toolbox. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 150–156, Online, July. Association for Computational Linguistics.
- Bañón, Marta, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online, July. Association for Computational Linguistics.
- Bojanowski, Piotr, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Dong, Daxiang, Hua Wu, Wei He, Dianhai Yu, and Haifeng Wang. 2015. Multi-task learning for multiple language translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1723–1732, Beijing, China, July. Association for Computational Linguistics.
- Escolano, Carlos, Marta R. Costa-jussà, José A. R. Fonollosa, and Mikel Artetxe. 2021. Multilingual machine translation: Closing the gap between shared and language-specific encoder-decoders. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 944–948, Online, April. Association for Computational Linguistics.
- Fan, Angela, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Çelebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. Beyond english-centric multilingual machine translation. *J. Mach. Learn. Res.*, 22:107:1–107:48.
- Goyal, Naman, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.
- Johnson, Melvin, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.
- Kingma, Diederik and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*, San Diego, CA, USA.
- Kocmi, Tom, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid).
- Kocmi, Tom, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet. In Koehn, Philipp, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore, December. Association for Computational Linguistics.
- Koehn, Philipp. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of*

- Machine Translation Summit X: Papers*, pages 79–86, Phuket, Thailand, September 13–15.
- Koppel, Kristina and Jelena Kallas. 2022. Eesti keele ühendkorpuste sari 2013–2021: mahukaim eesti-keelsete digitekstide kogu. *Eesti Rakenduslingvistika Ühingu aastaraamat = Estonian papers in applied linguistics*, 18:207–228.
- Korotkova, Elizaveta, Taïdo Purason, Agnes Luhtaru, and Mark Fishel. in press. Multilinguality or back-translation? A case study with Estonian. In *Accepted for publication at the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING)*. European Language Resources Association (ELRA).
- Kudo, Taku and John Richardson. 2018. Sentence-Piece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium, November. Association for Computational Linguistics.
- Kudugunta, Sneha, Isaac Caswell, Biao Zhang, Xavier Garcia, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. 2023. Madlad-400: A multilingual and document-level large audited dataset. In Oh, A., T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 67284–67296. Curran Associates, Inc.
- Lison, Pierre and Jörg Tiedemann. 2016. OpenSubtitles2016: Extracting large parallel corpora from movie and TV subtitles. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 923–929, Portorož, Slovenia, May. European Language Resources Association (ELRA).
- Lyu, Sungwon, Bokyung Son, Kichang Yang, and Jaekyoung Bae. 2020. Revisiting Modularized Multilingual NMT to Meet Industrial Demands. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5905–5918, Online, November. Association for Computational Linguistics.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loïc Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Searley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.
- Ott, Myle, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 48–53, Minneapolis, Minnesota, June. Association for Computational Linguistics.
- Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Post, Matt. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October. Association for Computational Linguistics.
- Purason, Taïdo, Andre Tättar, and Mark Fishel. 2024. Mixing and matching: Combining independently trained translation model components. In *Proceedings of the 1st Workshop on Modular and Open Multilingual NLP (MOOMIN 2024)*, pages 44–56, St Julians, Malta.
- Rei, Ricardo, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In Webber, Bonnie, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online, November. Association for Computational Linguistics.
- Rei, Ricardo, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In Koehn, Philipp, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Reimers, Nils and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation. In *Proceedings of the*

- 2020 *Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11.
- Schwenk, Holger, Vishrav Chaudhary, Shuo Sun, Hongyu Gong, and Francisco Guzmán. 2021a. WikiMatrix: Mining 135M parallel sentences in 1620 language pairs from Wikipedia. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1351–1361, Online, April. Association for Computational Linguistics.
- Schwenk, Holger, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela Fan. 2021b. CCMatrix: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6490–6500, Online, August. Association for Computational Linguistics.
- Sennrich, Rico, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 86–96, Berlin, Germany, August. Association for Computational Linguistics.
- Srivastava, Nitish, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958.
- Steinberger, Ralf, Bruno Poulliquen, Anna Widiger, Camelia Ignat, Tomaž Erjavec, Dan Tufiş, and Dániel Varga. 2006. The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC'06)*, Genoa, Italy, May. European Language Resources Association (ELRA).
- Tättar, Andre, Taido Purason, Hele-Andra Kuulmets, Agnes Luhtaru, Liisa Rätsep, Maali Tars, Mārcis Pinnis, Toms Bergmanis, and Mark Fishel. 2022. Open and competitive multilingual neural machine translation in production. In *Baltic Journal of Modern Computing*, volume 10, pages 422–434.
- Tiedemann, Jörg. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey, May. European Language Resources Association (ELRA).
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Ziemski, Michał, Marcin Junczys-Dowmunt, and Bruno Poulliquen. 2016. The United Nations parallel corpus v1.0. In Calzolari, Nicoletta, Khalid Choukri, Thierry Declerck, Sara Goggi, Marko Grobelnik, Bente Maegaard, Joseph Mariani, Helene Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 3530–3534, Portorož, Slovenia, May. European Language Resources Association (ELRA).

A Back-translated Dataset Sizes

The approximate number of sentence pairs in each of our back-translated corpora before filtering are shown in Table 7 (translated from Estonian into other languages) and Table 8 (translated from other languages into Estonian).

target language	ENC
Arabic	183.7
German	196.6
English	196.4
Spanish	172.7
Finnish	177.7
French	173.7
Lithuanian	174.0
Latvian	174.3
Russian	196.3
Swedish	171.4
Ukrainian	175.5
Chinese	189.0

Table 7: Sizes of the back-translation corpora translated from Estonian (unfiltered, in millions of sentence pairs). ENC stands for the Estonian National Corpus.

source language	corpus			
	NC	PC	UNPC	OS
Arabic	42.3	–	–	–
German	427.1	278.3	–	–
English	314.3	878.4	33.4	441.4
Spanish	72.1	208.4	–	–
Finnish	28.8	31.3	–	–
French	104.8	217.6	–	–
Lithuanian	7.6	13.2	–	–
Latvian	14.9	13.1	–	–
Russian	126.6	5.4	28.5	–
Swedish	–	49.1	–	–
Ukrainian	2.3	13.2	–	–
Chinese	13.9	14.2	–	–

Table 8: Sizes of the back-translation corpora translated into Estonian (unfiltered, in millions of sentence pairs). NC, PC, UNPC, and OS denote the NewsCrawl, ParaCrawl, United Nations Parallel Corpus, and OpenSubtitles corpora, respectively.

B Digital Object Identifiers for the Extended SynEst Corpus

The DOIs for each language pair of the extended SynEst corpus are shown in Table 9.

C Data Filtering

The back-translation datasets are filtered based on log probability of the generated translations. We only keep the examples that where log probability is higher than $\mu - 1.5\sigma$ where μ is the mean and σ is the standard deviation over all translation log probabilities for a given translation direction and corpus.

All data, both synthetic and parallel, are normalized with the MTee normalization script (Tättar et al., 2022) and filtered with OpusFilter (Aulamo et al., 2020). The following filters are used:

1. `LongWordFilter`: filter examples with words longer than 40 characters (default).
2. `LengthFilter`: filter examples longer than 1000 characters or shorter than 10 characters.
3. `LengthFilter`: filter examples longer than 100 words.
4. `LengthRatioFilter`: filter examples where the source and target sentence lengths differ more than 3 times in terms of number of words.
5. `CharacterScoreFilter` with threshold 1 (default) for the respective scripts.
6. `LanguageIDFilter` with `fastText` (Bojanowski et al., 2017) language identification model.
7. `LanguageIDFilter` with CLD2 language identification.
8. `TerminalPunctuationFilter` with the default parameters.
9. `NonZeroNumeralsFilter` with the default parameters.

This configuration is applied to all language pairs with the following exceptions:

- Arabic–Estonian, which uses filters 1 – 6 and uses minimal sentence length of 3 characters in filter 2;

language pair	DOI
Arabic–Estonian	doi.org/10.15155/y746-qa68
German–Estonian	doi.org/10.15155/2fy2-2k14
English–Estonian	doi.org/10.15155/5r1e-6r35
Spanish–Estonian	doi.org/10.15155/sqk9-ze70
Finnish–Estonian	doi.org/10.15155/hjw7-m565
French–Estonian	doi.org/10.15155/4vb6-ab11
Lithuanian–Estonian	doi.org/10.15155/7at2-jv07
Latvian–Estonian	doi.org/10.15155/erkh-k466
Russian–Estonian	doi.org/10.15155/4e20-vs27
Swedish–Estonian	doi.org/10.15155/jfws-ed89
Ukrainian–Estonian	doi.org/10.15155/xmpv-ft58
Chinese–Estonian	doi.org/10.15155/m6ww-j693
Estonian–all	doi.org/10.15155/ctz5-1d43

Table 9: DOIs for the extended SynEst corpus

- Chinese–Estonian, which only uses `LengthFilter` with maximal sentence length of 750 characters (no minimal length), `CharacterScoreFilter`, and `LanguageIDFilter` with `fastText` as language identification model.

Duplicates and test set overlaps are removed from the training dataset.

D Training Details

The models are trained with FairSeq (Ott et al., 2019). The NLLB-1.3B encoder consists of 24 transformer layers with embedding dimension 1024, feed-forward dimension 8192, and 16 attention heads. The decoders are randomly initialized and have 6 transformer layers; the dimensions of the decoders are the same as those of the encoder. The input and output embeddings of the decoder are shared. The vocabulary size is 256,000 for the encoder and 32,000 for the decoder (we train a separate SentencePiece (Kudo and Richardson, 2018) model for each output language). Models are trained on 8 GPUs (4 AMD MI250x 128GB GPU modules, each acting as 2 GPUs) with batch size 4,096 tokens per GPU. Models are trained for 2,000,000 updates, with checkpoints saved after every 2,000 updates. We use the inverse square root learning rate scheduler with 4,000 warm-up updates from initial learning rate 1×10^{-7} to maximum learning rate 5×10^{-4} . We use Adam optimizer (Kingma and Ba, 2015) with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Dropout probability (Srivastava et al., 2014) is 0.1, attention dropout 0.1, and activation

dropout is not used. The loss function is cross-entropy.

E COMET Scores

Tables 10, 11, and 12 show COMET scores (Rei et al., 2020) for translation from and into Estonian on the FLORES benchmark. Tables 13 and 14 contain results of translating the MTEE test sets from and into Estonian, respectively.

COMET scores were calculated with the default `wmt22-comet-da` model (Rei et al., 2022).

	target language					
	DE	EN	FI	RU	UK	ZH
NLLB-1.3B	84.19	88.33	86.65	86.71	85.90	80.01
MTEE	84.88	88.49	–	87.33	–	–
MADLAD-3B	84.64	<u>89.19</u>	89.03	85.55	82.23	<u>85.57</u>
Ours	<u>85.95</u>	88.92	<u>90.25</u>	<u>88.26</u>	<u>87.79</u>	84.51
DeepL	87.08	89.54	91.44	89.67	90.07	87.69
Google	87.21	89.75	90.70	89.74	89.77	87.78

Table 10: COMET scores on the FLORES-devtest benchmark for models translating from Estonian into other languages. The best scores overall are shown in **bold**, and the best scores among open-source models are underlined. For our models, we report the score of the checkpoint with the best validation BLEU. With MTEE, we use the general-domain model to translate the FLORES benchmark.

	source language					
	AR	DE	EN	ES	FI	FR
NLLB-1.3B	84.08	87.37	89.36	86.13	87.23	87.00
MTEE	–	88.82	89.34	–	–	–
MADLAD-3B	<u>87.65</u>	88.86	90.65	87.78	88.84	88.01
Ours	87.34	<u>90.42</u>	<u>91.60</u>	<u>88.67</u>	<u>90.58</u>	<u>89.76</u>
DeepL	89.02	91.25	92.54	89.78	91.13	90.67
Google	88.35	90.34	91.77	89.29	90.72	90.17

Table 11: COMET scores on the FLORES-devtest benchmark for models translating from Arabic, German, English, Spanish, Finnish, and French into Estonian. The best scores overall are shown in **bold**, and the best scores among open-source models are underlined. For our model, we report the score of the checkpoint with the best validation loss (the same checkpoint is used for all source languages). With MTEE, we use the general-domain model to translate the FLORES benchmark.

	source language					
	LT	LV	RU	SV	UK	ZH
NLLB-1.3B	85.36	85.78	86.27	87.50	85.69	84.03
MTEE	–	–	88.28	–	–	–
MADLAD-3B	87.82	<u>90.27</u>	86.07	88.54	83.44	<u>88.48</u>
Ours	<u>88.72</u>	89.92	<u>89.37</u>	<u>90.57</u>	<u>89.08</u>	88.18
DeepL	90.23	91.05	89.92	91.55	90.15	89.91
Google	89.68	90.46	89.42	90.77	89.24	89.55

Table 12: COMET scores on the FLORES-devtest benchmark for models translating from Lithuanian, Latvian, Russian, Swedish, Ukrainian, and Chinese into Estonian. The best scores overall are shown in **bold**, and the best scores among open-source models are underlined. For our model, we report the score of the checkpoint with the best validation loss (the same checkpoint is used for all source languages). With MTEE, we use the general-domain model to translate the FLORES benchmark.

	ET-DE	ET-EN	ET-RU
News			
NLLB-1.3B	83.35	83.64	85.15
MTEE	85.12	84.03	86.70
MADLAD-3B	83.74	<u>85.07</u>	82.41
Ours	<u>85.41</u>	84.19	<u>87.36</u>
DeepL	86.32	85.51	88.71
Google	86.64	85.25	88.70
Crisis			
NLLB-1.3B	83.79	85.08	87.60
MTEE	<u>85.62</u>	86.76	90.18
MADLAD-3B	81.26	<u>87.00</u>	85.75
Ours	85.50	86.65	90.77
DeepL	86.18	87.88	90.62
Google	86.26	88.39	90.24
Military			
NLLB-1.3B	83.05	86.35	88.72
MTEE	84.26	87.14	89.88
MADLAD-3B	80.62	<u>87.34</u>	85.85
Ours	85.54	87.04	<u>90.53</u>
DeepL	84.68	87.51	90.51
Google	85.12	88.12	90.60
Legal			
NLLB-1.3B	84.51	87.12	90.84
MTEE	86.72	88.17	92.33
MADLAD-3B	85.01	88.01	90.85
Ours	<u>87.04</u>	88.14	92.39
DeepL	87.09	87.91	91.07
Google	86.68	87.62	91.32
Spoken			
NLLB-1.3B	80.55	81.65	83.30
MTEE	82.22	82.19	84.04
MADLAD-3B	81.85	<u>83.75</u>	81.44
Ours	<u>82.92</u>	81.96	<u>84.37</u>
DeepL	83.21	83.94	86.08
Google	83.98	84.26	85.67

Table 13: COMET scores on the MTEE domain benchmark sets for models translating from Estonian into other languages. The best scores overall are shown in **bold**, and the best scores among open-source models are underlined. For our models, we report the score of the checkpoint with the best validation BLEU. With MTEE, we calculate the scores on the same model outputs as used by Tättar et al. (2022).

	DE-ET	EN-ET	RU-ET
News			
NLLB-1.3B	85.80	86.61	87.41
MTEE	87.83	85.85	89.34
MADLAD-3B	87.32	87.96	87.00
Ours	<u>89.88</u>	88.93	91.07
DeepL	90.45	89.93	90.53
Google	90.00	88.47	90.36
Crisis			
NLLB-1.3B	89.55	91.08	87.75
MTEE	91.00	93.96	91.91
MADLAD-3B	91.48	93.02	88.86
Ours	93.83	94.51	92.44
DeepL	92.52	94.36	91.81
Google	92.25	94.07	91.49
Military			
NLLB-1.3B	88.73	92.26	89.24
MTEE	90.81	93.40	92.00
MADLAD-3B	90.33	92.92	89.16
Ours	92.56	<u>93.55</u>	92.57
DeepL	92.19	94.28	91.81
Google	91.43	93.92	91.54
Legal			
NLLB-1.3B	90.07	92.88	91.13
MTEE	91.96	95.50	94.23
MADLAD-3B	92.84	93.50	93.54
Ours	94.51	95.62	94.72
DeepL	93.49	95.45	93.49
Google	92.35	94.54	92.22
Spoken			
NLLB-1.3B	86.59	88.31	84.26
MTEE	89.56	90.15	87.51
MADLAD-3B	89.11	90.13	84.33
Ours	90.88	<u>90.75</u>	88.47
DeepL	90.06	90.98	87.23
Google	89.58	90.72	87.30

Table 14: COMET scores on the MTEE domain benchmark sets for models translating from other languages into Estonian. The best scores overall are shown in **bold**, and the best scores among open-source models are underlined. For our model, we report the score of the checkpoint with the best validation loss (the same checkpoint is used for all source languages). With MTEE, we calculate the scores on the same model outputs as used by Tättar et al. (2022).

The EAMT organisers gratefully acknowledge the support from the following sponsors.

Silver



Bronze



Collaborators



Supporters

SPRINGER NATURE

Media Sponsors



Author Index

- Aditya Pavani, Penumalla, 207
Al-Ramadan, Ali, 480
Alvarez-Vidal, Sergi, 573
Appicharla, Ramakrishna, 246
Aranberri, Nora, 548
Aymo, Mahmoud, 373
- Ballier, Nicolas, 396
Bane, Fred, 373
Batista-Navarro, Riza, 590
Bawden, Rachel, 431
Bayliss, Chris, 561
Bechara, Hannah Dorothy, 623
Bentivogli, Luisa, 300
Berger, Nathaniel, 636
Bhaskar, Yash, 207
Bhattacharyya, Pushpak, 246
Blanch Miró, Tània Blanch Miró, 373
Bouillon, Pierrette, 387
Break, Page, 1, 8, 315, 560
Briva-Iglesias, Vicent, 444
Brockmann, Justus, 455
Buschbeck, Bianka, 610
Bénard, Maud, 431
- Cadwell, Patrick, 347
Camargo, João Lucas Cavalheiro, 492
Canavese, Paolo, 347
Carl, Michael, 480
Castilho, Sheila, 492
Charkiewicz, Adrian, 80
Chen, Pinzhen, 181
Chen, Yongjian, 229
Chereji, Raluca, 455
Chiocchetti, Elena, 573
Chiruzzo, Luis, 116
Chu, Chenhui, 147
Ciobanu, Dragoş, 455
Clergerie, Éric Villemonte De La, 431
Costa-jussà, Marta R., 37, 360
- Dabre, Raj, 147
Dale, David, 360
Deguchi, Hiroyuki, 191
Deilen, Silvana, 469
Diego, Sara, 580
Dinarelli, Marco, 396
Dinh, Tu Anh, 133
- Dutta Chowdhury, Koel, 411
- Einarsson, Hafsteinn, 24
Ekbal, Asif, 246
Elbayad, Maha, 360
Erofeev, Gleb, 337
Esamotunu, Raphaël, 431
Eschbach-Dymanus, Johannes, 610
Escolano, Carlos, 37
España-Bonet, Cristina, 411
Esperança-Rodier, Emmanuelle, 396
Essenberger, Frank, 610
Etchegoyhen, Thierry, 9
Exel, Miriam, 610, 636
- Fantinuoli, Claudio, 327
Farinhas, António, 258
Farrús, Mireia, 229
Fernandes, Patrick, 258
Fishel, Mark, 647
Friberg, Théo, 100
- Gaido, Marco, 2
Gain, Baban, 246
García Gilabert, Javier, 37
Garrido, Sergio Hernandez, 469
Genabith, Josef Van, 411
Gete, Harritxu, 9
Gladkoff, Serge, 337
Gonzalez-Saez, Gabriela, 396
Gromann, Dagmar, 507
Guo, Zhicheng, 181
Guttman, Kamil, 80
- Haddow, Barry, 181
Han, Lifeng, 337, 590
Hauhio, Iikka, 100
He, Jianfei, 68
He, Sui, 316, 396
Heafield, Kenneth, 181
Hiebl, Bettina, 507
Hong, Kung Yin, 590
Huck, Matthias, 636
Huguin, Mathilde, 431
Hörner, Julian, 469
- Iranzo-Sánchez, Javier, 4

Jankin, Slava, 623
 Jia, Xiaohua, 68

 King, Adam, 59
 Korotkova, Elizaveta, 647
 Krishnamurthy, Parameswari, 207
 Kunilovskaya, Maria, 411
 Kurohashi, Sadao, 147
 Kübler, Natalie, 431

 Lagzdiņš, Andis, 600
 Lai, Huiyuan, 286
 Lapshinova-Koltunski, Ekaterina, 469
 Li, Wenjie, 68
 Li, Zhijian, 229
 Lopez, Fabien, 396
 Lu, Sheng, 480

 Maaß, Christiane, 469
 Manohara, Krishnamoorthy, 623
 Martins, Andre, 258
 McGill, Euan, 116
 Mestivier, Alexandra, 431
 Meyer, Francois, 147
 Michelot, Mona, 431
 Moniz, Helena Silva, 628
 Moorkens, Joss, 492
 Mujadia, Vandan, 207

 Nagata, Masaaki, 191
 Nakhle, Mariam, 396
 Negri, Matteo, 300
 Nenadic, Goran, 337, 590
 Niehues, Jan, 133
 Nowakowski, Artur, 80, 164
 Nunziatini, Mara, 580

 Oakley, Chris, 561
 Ohuoba, Adaeze Ngozi, 537
 Oliver, Antoni, 573
 O'Brien, Sharon, 444

 Pal, Santanu, 246
 Palzer, Tobias, 133
 Peng, Ziqian, 431
 Piergentili, Andrea, 300
 Piperidis, Stelios, 275
 Ploeger, Esther, 286
 Pokrywka, Mikołaj, 80
 Prescott, Charlotte, 561
 Przybyl, Heike, 411

 Qader, Raheel, 396

 Ramos, Miguel Moura, 258
 Riezler, Stefan, 636
 Rios, Miguel, 455
 Romary, Laurent, 431
 Rossi, Caroline, 396
 Rostek, Zofia, 164
 Roussis, Dimitris, 275

 Saggion, Horacio, 116
 Savoldi, Beatrice, 300
 Scarton, Carolina, 561
 Schwab, Didier, 396
 Secară, Alina, 455
 Sharma, Dipti, 207
 Sharoff, Serge, 537
 Shravya, Kukkapalli, 207
 Silveira, Diana, 628
 Simonsen, Annika, 24
 Sofianopoulos, Sokratis, 275
 Song, Haiyue, 147
 Sorokina, Irina, 337
 Stemle, Egon, 573
 Sun, Shichao, 68

 Tanaka, Hideki, 147
 Theel, Vanessa, 469
 Toral, Antonio, 229, 286
 Torres, João, 373
 Torrón, Marina Sánchez, 628
 Turner, James Robert, 396

 Uguet, Celia Soler, 373
 Urlana, Ashok, 207

 Vacheva, Neil, 600
 Vamvas, Jannis, 6
 Van Noord, Rik, 286
 Vasiljevs, Andrejs, 600
 Vincent, Sebastian, 561
 Volkart, Lise, 387
 Vīksna, Rinalds, 600

 Walker, Callum, 537
 Wang, Xiaoman, 327
 Watanabe, Taro, 191
 Wiesinger, Claudia, 455
 Wisniewski, Dawid, 164

Yang, Jun, 396
YU, Bokai, 360
Yvon, François, 431

Zhu, Lichao, 431
Ziemer, Sophie, 469

Zaretskaya, Anna, 373