# LREC-COLING 2024

**The Fourth Workshop on
Human Evaluation of NLP Systems
(HumEval 2024)**

Workshop Proceedings

Editors

Simone Balloccu, Anya Belz, Rudali Huidrom, Ehud Reiter,
João Sedoc and Craig Thomson

21 May, 2024
Torino, Italia

**Proceedings of the Fourth Workshop on Human Evaluation of NLP Systems (HumEval 2024)**

Jointly organized by the ELRA Language Resources Association
and the International Committee on Computational Linguistics

# Preface

Welcome to HumEval 2024!

We are pleased to present the proceedings of the fourth workshop on Human Evaluation of NLP Systems (HumEval) which is taking place as part of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024).

Human evaluation is vital in NLP, and it is often considered as the most reliable form of evaluation. It ranges from large-scale crowd-sourced evaluations to the much smaller experiments routinely encountered in conference papers. With this workshop we wish to create a forum for current human evaluation research, a space for researchers working with human evaluations to exchange ideas and begin to address the issues that human evaluation in NLP currently faces, including aspects of experimental design, reporting standards, meta-evaluation and reproducibility.

We are truly grateful to the authors of the submitted papers that showed interest in human evaluation research. The HumEval workshop accepted 8 submissions. The accepted papers cover a broad range of NLP areas where human evaluation is used: machine translation, natural language generation, text simplification, conversational search. Several papers are addressing reproducibility of human evaluations. The workshop once again hosted the results session of the ReproNLP Shared Task on Reproducibility of Evaluations in NLP which consisted of the presentation of overall results by the organizers and 18 oral and poster presentations by participants.

This workshop would not have been possible without the hard work of the program committee. We would like to express our gratitude to them for writing detailed and thoughtful reviews in a very constrained span of time. We also thank our invited speakers, Mark Diaz and Sheila Castilho, for their contribution to our program. We are grateful for the help from the LREC-COLING 2024 workshop organizers, and to all the people involved in setting up the infrastructure.

You can find more details about the worskhop on its website: `https://humeval.github. io/`.
Simone, Craig, João, Anya, Ehud, Rudali

# Organizing Committee

**Organizers**

Simone Balloccu, ADAPT Centre, Dublin City University, Ireland
Anya Belz, ADAPT Centre, Dublin City University, Ireland
Rudali Huidrom, ADAPT Centre, Dublin City University, Ireland
Ehud Reiter, University of Aberdeen, UK
João Sedoc, New York University, US
Craig Thomson, University of Aberdeen, UK

**Program Committee**

Gavin Abercrombie, Heriot-Watt University, UK
Jose Maria Alonso-Moral, University of Santiago de Compostela, ES
Mohammad Arvan, University of Illinois, Chicago, US
Anouck Braggaar, Tilburg University, NL
Daniel Braun, University of Twente, NL
Javier Corbelle, University of Santiago de Compostela, ES
Tanvi Dinkar, Heriot-Watt University, UK
Ondrej Dusek, Karls University, CZ
Steffen Eger, University of Mannheim, DE
Manuela Hürlimann, Zurich University of Applied Sciences, CH
Mateusz Lango, Poznań University of Technology, PL
Yiru Li, Groningen University, NL
Michela Lorandi, Dublin City University / ADAPT, IE
Saad Mahamood, trivago, DE
Zola Mahlaza, University of Cape Town, ZA
Gonzalo Mendez, University of Madrid, ES
Margot Mieskes, University of Applied Sciences, Darmstadt, DE
Jie Ruan, Peking University, CN
Patricia Schmidtova, Karls University, CZ
Raj Shah, Georgia Tech, US
Barkavi Sundararajan, University of Aberdeen, UK
Supryadi, Tianjin University, CN
Chris van der Lee, Tilburg University, NL Xiaojun Wan, Peking University, CN
Deyi Xiong, Tianjin University, CN
Emiel van Miltenburg, Tilburg University, NL
Chuang Liu, Tianjin University, CN

**Invited Speakers**

Mark Diaz, Google Research
Sheila Castilho, ADAPT/DCU

# Table of Contents

# Conference Program

**May 21, 2024 (continued)**

| | |
|---|---|
| 14:00–14:45 | **Oral Session 2** |

*Insights of a Usability Study for KBQA Interactive Semantic Parsing: Generation Yields Benefits over Templates but External Validity Remains Challenging*
Ashley Lewis, Lingbo Mo, Marie-Catherine de Marneffe, Huan Sun and Michael White

*Extrinsic evaluation of question generation methods with user journey logs*
Elie Antoine, Eléonore Besnehard, Frederic Bechet, Geraldine Damnati, Eric Kergosien and Arnaud Laborderie

*Towards Holistic Human Evaluation of Automatic Text Simplification*
Luisa Carrer, Andreas Säuberli, Martin Kappus and Sarah Ebling

| | |
|---|---|
| 14:45–16:00 | **ReproNLP Shared Task Session 2** |

| | |
|---|---|
| 16:00–16:30 | **Coffee Break** |

| | |
|---|---|
| 16:30–17:15 | **Invited Talk 2** |

*All That Agrees Is Not Gold: Evaluating Ground Truth and Conversational Safety*
Mark Diaz

| | |
|---|---|
| 17:15–18:00 | **Oral Session 3** |

*Decoding the Metrics Maze: Navigating the Landscape of Conversational Question Answering System Evaluation in Procedural Tasks*
Alexander Frummet and David Elsweiler

| | |
|---|---|
| 18:00–18:05 | **Closing Remarks** |

# Quality and Quantity of Machine Translation References for Automatic Metrics

**Vilém Zouhar**
ETH Zürich
vzouhar@ethz.ch

**Ondřej Bojar**
Charles University
bojar@ufal.cuni.cz

### Abstract

Automatic machine translation metrics typically rely on *human* translations to determine the quality of *system* translations. Common wisdom in the field dictates that the human references should be of very high quality. However, there are no cost-benefit analyses that could be used to guide practitioners who plan to collect references for machine translation evaluation. We find that higher-quality references lead to better metric correlations with humans at the segment-level. Having up to 7 references per segment and taking their average (or maximum) helps all metrics. Interestingly, the references from vendors of different qualities can be mixed together and improve metric success. Higher quality references, however, cost more to create and we frame this as an optimization problem: given a specific budget, what references should be collected to maximize metric success. These findings can be used by evaluators of shared tasks when references need to be created under a certain budget.

## 1. Introduction

Machine translation systems are robustly evaluated through human annotation. This is non-scaleable and non-replicable (Freitag et al., 2021a) for settings such as shared tasks where a number of teams submit automatic translations of the same testset. Automatic metrics aim to provide a cheap and replicable solution. Given the translation and possibly the source and reference segments, they produce a score that correlates with what a human annotator would predict. There is support and evidence for not using references (Lommel, 2016) in metrics, i.e. quality estimation (Specia and Shah, 2018; Rei et al., 2021). Still, most of the commonly used metrics (Section 3) require human reference translations (Freitag et al., 2023). These metrics work by comparing either the overlap on the surface-level (e.g. BLEU, Papineni et al., 2002), of semantic representations (e.g. COMET, Rei et al., 2020) or some downstream task (e.g. MTEQA, Krubiński et al., 2021).

Humans also do not always arrive at perfect translations and thus the quality of the references themselves varies (Castilho et al., 2018). In cases of very poor translations, such as non-translation,[1] the reference-based metrics would clearly fail. While low-quality references are known to decrease the metric correlations (Freitag et al., 2023), the extent of this effect and interactions with other phenomena remains unclear. Many automatic machine translation metrics support multiple references for a single translation natively or by using an aggregation such as the average. For phrase-based MT and BLEU, the trade-off between the number of references vs. the test set size was studied by

Bojar et al. (2013, Section 5), concluding that a single-reference test set of 3000 sentences can be comparable to 6–7 references with just 100–200 test sentences. The usefulness of multiple references was later disputed (Freitag et al., 2020) for state-of-the-art system evaluation and some recent metrics do not even support multiple references. Additionally, a professional experienced translator is likely to produce a better translation than an average crowd-worker. However, the cost of a high-quality human translation is likely also much higher.

In this paper, we aim to quantify the trade-off between reference **quality, quantity** and **cost** for segment-level automatic metric performance. We base our experiments on a small-scale English→Czech dataset with multiple references of varying qualities.

We pose research questions with immediate implications for practitioners. The short answers here are only summaries.

**Q:** Are higher-quality references useful for automatic evaluation?
**A:** Low-quality should be avoided. Too much investment has diminishing returns. (Sec. 4.1)
**Q:** Are multiple references useful?
**A:** Yes. Averaging or taking the maximum across reference improves the metrics. (Sec. 4.2)
**Q:** How to allocate the budget?
**A:** By not focusing exclusively on either quality or quantity of references but their combination. This can be computed by Algorithm 1, given a list of vendors and their attributes. (Sec. 4.4)

## 2. Related Work

Reference quality is known to affect machine translation evaluation. Freitag et al. (2023) note that very low-quality references reduce metric success.

---

[0]github.com/ufal/optimal-reference-translations
hf.co/datasets/zouharvi/optimal-reference-translations
[1]Text left untouched in the source language.

This stands in contrast to the pre-neural machine translation era where the reference quality did not play an important role in certain settings (Hamon and Mostefa, 2008). This is likely caused by the much higher quality of systems being compared. Vernikos et al. (2022) hypothesize that ambiguous and vague references are the culprits of metric success deterioration. Additionally, Freitag et al. (2020) study how to avoid low-quality references in human translation campaigns.

The BLEU metric (Papineni et al., 2002) was intended to be used with multiple metrics, which was only rarely put in practice over the years. Nevertheless, newer and more sophisticated methods exist to incorporate them (Qin and Specia, 2015). Our results from Figure 1 confirm the older observations of Finch et al. (2004) or Bojar et al. (2013, Section 5) who study the effect of the reference count on metric performance. Finally, multiple references can be used in training better machine translation systems (Madnani et al., 2008; Zheng et al., 2018; Khayrallah et al., 2020; Mi et al., 2020) or for analyzing model uncertainty (Ott et al., 2018) or evaluation uncertainty (Zhang and Vogel, 2004, 2010; Fomicheva et al., 2020). It is also used outside of machine translation for measuring consensus (Vedantam et al., 2015).

The budget allocation algorithm is reminiscent of active learning or data selection. In machine translation, this is limited to selecting training examples (Haffari et al., 2009; González-Rubio et al., 2012; van der Wees et al., 2017; Shi and Huang, 2020; Mendonça et al., 2023). We focus on algorithmic data selection for higher-quality *evaluation*. We aim to complete similar works on practical advice on machine translation evaluation. Kocmi et al. (2021, 2024) study the reliability of metrics from the perspective of deployment decisions. We show that the configuration of references can make these metrics stronger or weaker on segment-level.

## 3. Setup

To evaluate the effect of references on automatic machine translation evaluation, we need data with controlled references and reference-based metrics.

**Optimal Reference Translations.** Zouhar et al. (2023) re-annotate a subset of the English→Czech testset from the News domain of the WMT2020 campaign (Barrault et al., 2020). New references were created by translating the original source in 4 different human settings ranging from generic translation vendors to translatology academics following a novel protocol leading to so-called "optimal reference translations" (Kloudová et al., 2023). This phase was followed by a human annotation and post-editing phase performed by 11 annotators of varying professionalities.

Zouhar et al. (2023) study whether the human quality of the references is really the highest achievable one. They stop short of evaluating the impact of this on machine translation evaluation. We re-purpose their data and system submissions from Barrault et al. (2020). We refer to the references, from lowest to highest quality of the source, as R1, R2, R3, and R4. Specifically, R1 to R2 come from standard translation vendors,[2] R3 is high-quality translation vendor, and R4 is the work of translatologists (the optimal reference). See Table 1 for basic statistics.

| Source segments & documents | 160 & 20 |
|---|---|
| Average source segment length | 34 tokens |
| Reference segments | $160 \times 4 = 640$ |
| Reference post-editing | $160 \times 4 \times 11 = 7040$ |
| Systems & system segments | 13 & $160 \times 13 = 2080$ |

Table 1: Overview of the used dataset.

**Automated Metrics.**[3] For the metrics, we use BLEU (Papineni et al., 2002), chrF (Popović, 2015), TER (Snover et al., 2006), COMET$_{20}$ (Rei et al., 2020), its referenceless version COMET$_{20}^{QE}$, and its updated iteration COMET$_{22}$ (Rei et al., 2022), and BLEURT (Sellam et al., 2020). We select these as a representative set of widely-used string-matching and trainable metrics.

**Metric Evaluation.** We focus on and evaluate metric success at the segment-level ("sentence"-level) by correlating the metric scores with human scores using Kendall's $\tau$.[4] Each translation receives a human score and automatic metric scores which are correlated. This is the standard segment-level evaluation adopted by the WMT Metrics Shared Task (Freitag et al., 2021b, 2022, 2023). In our case (WMT2020), the human segment-level judgments were created from Direct Assessment (Graham et al., 2016) judgements following the "DARR" conversion as described by Mathur et al. (2020): Candidate translations from MT systems were scored on their own, independently of other candidates. For each pair of judgements of candidates translating the same source, we construct one golden-truth item of pairwise comparison if the two individual scores differ by more than 25% absolute. As Mathur et al. (2020), we believe that this difference in the judgement is big enough to trust the simulated pairwise comparison.

---

[2]Nevertheless, based on observations of Kloudová et al. (2021), R1 are to a large extent post-edits of one of the participating systems.

[2] `BLEU|#:1|c:mixed|e:yes|tok:13a|s:exp`
`chrF|#:1|c:mixed|e:yes|nc:6|nw:0|s:no`
`TER|#:1|c:lc|t:tercom|nr:no|pn:yes|as:no`

[4]$\tau = $ (#concordant − #discordant)/#pairs; read more on the definition in Macháček and Bojar (2014).

| Metric | R1 | R2 | R3 | R4 |
|---|---|---|---|---|
| BLEU | 0.082 ▁ | 0.103 ▁ | 0.109 ▁ | 0.103 ▁ |
| chrF | 0.090 ▁ | 0.125 ■ | 0.128 ■ | 0.123 ■ |
| TER | 0.082 ▁ | 0.092 ▁ | 0.114 ■ | 0.105 ■ |
| $COMET_{20}$ | 0.172 ■ | 0.176 ■ | 0.185 ■ | 0.181 ■ |
| $COMET_{22}$ | 0.189 ■ | 0.195 ■ | 0.191 ■ | 0.192 ■ |
| BLEURT | 0.159 ■ | 0.156 ■ | 0.199 ■ | 0.178 ■ |
| **Average** | 0.129 ■ | 0.141 ■ | 0.154 ■ | 0.147 ■ |
| $COMET_{20}^{QE}$ | | 0.171 ■ | | |

Table 2: Segment-level Kendall's $\tau$ between automatic metrics and human scores. The metrics are computed with respect to each of the four references. The black boxes indicate the value visually and comparable across both columns and rows. ◉ The R3 translation yields the best results as the reference, despite not being the optimal translation from the human perspective.

| Proficiency | $R1^{PE}$-R1 | $R2^{PE}$-R2 | $R3^{PE}$-R3 | $R4^{PE}$-R4 |
|---|---|---|---|---|
| **Layman** | +0.019 | +0.011 | +0.011 | +0.011 |
| **Student** | +0.009 | +0.005 | +0.001 | −0.002 |
| **Professional** | +0.025 | +0.011 | +0.004 | +0.002 |

Table 3: Difference in Kendall's $\tau$ between using original translations (in Table 2) and their post-edited versions. The post-editing comes from translators on different levels. The correlations are averaged across all metrics; see Tables 6 and 10 for per-metric breakdowns. ◉ In most cases, using post-edited versions improves metric performance.

# 4. Experiments

## 4.1. Reference Quality is Important

As stated in Section 3, we have access to four human translations of varying quality. In Table 2, we show the metric success measured by correlation with human scores. The metrics stay the same but the references they use are changed. Across both string-matching and parametrized model-based metrics, R1, the worst human translation, leads to the worst metric performance. The best performance is achieved with R3, a standard professional translation. Notably, it is not R4 which was created by professional translatologists and was also the most expensive one. This can be explained by the presence of translation shifts, which occur more frequently on this professionality level, but can negatively impact the utility of the reference (Fomicheva, 2017). Translation shifts in general refer to deviation from the original structure or meaning. For our new references, the translatologists paid attention to preserve the meaning but they often restructured the sentences. They did this to avoid translationese as much as possible and to express the subtleties of information structure (given

| Aggregation | R3 | R{3,4} | Rx | $Rx^{PE}$ |
|---|---|---|---|---|
| **Average** | 0.154 ■ | 0.159 ■ | 0.166 ■ | 0.164 ■ |
| **Max** | 0.154 ■ | 0.155 ■ | 0.165 ■ | 0.167 ■ |

Table 4: Average performance of metrics with multiple references. See Table 9 for per-metric breakdown. ◉ All aggregation methods improve the performance over the best single one, R3.

vs. new information) which is natively expressed via Czech word order. These boosted word order differences make it harder for automatic metrics to match the, rather translationese, candidate and the reference. We anticipate that more fluent large language model-based MT could sound less translationese and the "optimal reference translations" will serve better in this setting. See Section 5 for an example and analysis.

A simple way of improving a translation is to post-edit (refine) it, which is cheaper that translating it from scratch (Daems and Macken, 2020; Zouhar et al., 2021). Moreover, Bojar et al. (2013, Figure 7) show that such post-edited references lead to a better performance of BLEU, because "every n-gram mismatch indicates an error". With standard references, an n-gram mismatch often means just lack of reference coverage. However, such post-edited references need to be ideally created for each evaluated MT system. In our case, the post-edits were created starting from *human* reference translations Rx. We mark them $Rx^{PE}$ and use them as references for the automatic metrics in Table 3. The post-editors are either laymen with knowledge of both languages, students of translatology, or professional translators. While the proficiency level plays a role, in most cases the post-edited translations serve as better references. Table 6 below lists the raw metric score changes in a closer detail.

## 4.2. Multiple References are Useful

The previous section provided an analysis of how individual references affect metric performance. In many situations, however, multiple references are available. While some metrics, such as BLEU, support multiple references natively, one can also aggregate them using either segment-level averages or maxima (i.e. compute multiple scores for each segment and take the average or maximum). In Table 4 we consider three setups: (1) two high-quality references, R3 and R4, (2) all human translations, Rx, or (3) all post-edited human translations, $Rx^{PE}$. Across all metrics, this segment-level aggregation improves the correlation with humans, especially in the case of using the original four human translations. Taking the maximum and not the average has the advantage that there exists a specific reference which yields that particular score. The

Figure 1: Metric performance with multiple sampled references from the pool of the original human translations and their post-edited versions. Confidence t-test intervals indicate 99% confidence of the mean (of 10 samples) being in the shaded area. ⊙ Biggest advantage is gained from at least three references and taking their segment-level average (max aggregation not shown).

maximum also reflects the spirit of automated evaluation: measuring some similarity between the candidate and reference translations. With more references, taking the maximum corresponds to first finding the most similar reference. We include all subsets as references in Table 9.

To systematically study the effect of reference count on metric performance, in Figure 1 we randomly sample $x$ references from the whole pool of original and post-edited translations, irrespective of their quality. The biggest gains in metric performance are achieved until seven references and further gains are negligible, which is in line with the observations of Bojar et al. (2013, Section 5).

| Metric | R1 | R2 | R4 | R3 | R1$^{PE}$ | R3$^{PE}$ |
|---|---|---|---|---|---|---|
| BLEU | 24.2 | 31.5 | 27.3 | 37.1 | 23.9 | 31.0 |
| chrF | 55.7 | 60.3 | 56.1 | 63.0 | 54.5 | 58.4 |
| TER | -63.3 | -53.0 | -59.4 | -48.7 | -64.1 | -58.9 |
| COMET$^{20}$ | 65.5 | 68.9 | 61.0 | 68.2 | 60.4 | 61.4 |
| COMET$^{22}$ | 84.6 | 84.9 | 83.6 | 84.8 | 83.6 | 83.7 |
| BLEURT | 61.3 | 66.1 | 64.5 | 68.8 | 61.6 | 64.9 |

Table 5: Raw average scores across metrics and references. TER scores are flipped to make higher numbers be better. The columns are sorted by quality of references from worst to best as reported in Table 2. ⊙ For most metrics, higher absolute metric scores correspond to better evaluation (numbers are growing from left to right), except for post-edited human references Rx$^{PE}$ which serve better as references (are more to the right) but lead to lower absolute metric scores.

## 4.3. References and Metric Scores

To understand the effect of different metrics, we show the average *raw* scores of each metric in Table 5. While it appears that the higher the raw

score, the better the metric performance (low score of R1 and high scores of R3 and R4), this trend does not explain the improvements of using the post-edited versions, e.g. as R1$^{PE}$ over R1, or R3$^{PE}$ over R3. In fact, the post-edited versions always lead to lower raw scores. This could be the result of either further translation shifts as the post-edits are based on a translation and not the source or additional (fully justified) corrections in the references which lead to fewer matches with the candidates.

## 4.4. Allocating a Budget for References

Usually, it is simple to gather many source sentences and let multiple systems translate them. Evaluating all of them using human annotators is unattainable but running automatic metrics is not. However, these require references, which are also costly. It remains unclear how many references and of which quality to obtain to achieve the most reliable automatic quality assessment under a given budget.



Figure 2: Metric performance with references (ordered by usefulness) from mixed sources (e.g. 25% R1 and 75% R2; rightmost is 100% R3). ⊙ Mixing references does not hurt any metric.

**Can references be mixed?** To assess what types of configurations of references *can* lead to the most reliable automatic evaluation, we first validate if references can be meaningfully mixed. For example, if it is viable that 75% of the sources can have references from a cheaper vendor R1 and 25% from a higher-quality but more expensive vendor R3. This is different from Table 4 where each segment had exactly two references from the same two sources. In Section 4.3, we show that using lower-quality references R1 leads to lower absolute metric scores (e.g. BLEU = 24.2) as opposed to higher-quality ones R3 (e.g. BLEU = 37.1). This holds across all metrics. Bojar et al. (2010) observe that lower BLEU scores are less reliable, but they refer to the range of BLEU < 20. It is thus questionable if BLEUs at 20–40 correlate differently with human MT quality judgements. In Figure 2, we mix some of the references together for evaluation, but staying at single-reference evaluation.

**Algorithm 1** Budget Allocation for References

**Input:** Source segments $S$, levels $L$, cost function $\textsc{Cost} : L \to \mathbb{R}^+$, utility function
$\quad\quad\textsc{Util.} : L \to \mathbb{R}^+$, tradeoff hyperparameter $\lambda \in [0,1]$, temperature $t > 0$, budget $B \in \mathbb{R}^+$.
**Output:** Assignment $R : L \to 2^S$.
**Note:** Figure 4 contains a patience mechanism instead of exit on error.

1: $L \leftarrow \textsc{Sort}(L, \textsc{Cost})$
2: $R[L_0] \leftarrow S;\quad\quad O \leftarrow R$ ▷ Assign everything to the cheapest level at first.
3: **while** $\sum_{l \in L} |R[l]| \cdot \textsc{Cost}(l) < B \;\wedge\;$ **no exception do**
4: $\quad O \leftarrow R$
5: $\quad a \sim \textsc{Sample}(\textsc{Promote} : \lambda, \textsc{Add} : 1 - \lambda)$ ▷ Select action.
6: $\quad X^+ \leftarrow \{\langle s, l \rangle | l \in L, s \in S \setminus R[l]\}$ ▷ Samples that could be added to $R[l]$.
7: $\quad X^- \leftarrow \{\langle s, l \rangle | l \in L, s \in R[l]\}$ ▷ Samples that could be removed from $R[l]$.

8: $\quad$ **if** $a = \textsc{Add}$ **then**
9: $\quad\quad x, l \sim \textsc{Sample}(\{\langle x, l \rangle : \frac{\sigma(\textsc{Util.}(l) - \textsc{Cost}(l))^{1/t}}{Z} | x, l \in X^+\})$ ▷ Sample a segment to add.
10: $\quad\quad R[l] \leftarrow R[l] \cup \{x\}$ ▷ Commit transaction.

11: $\quad$ **else if** $a = \textsc{Promote}$ **then**

12: $\quad\quad x^+, l^+ \sim \textsc{Sample}(\{\langle x, l \rangle : \frac{\sigma(\textsc{Util.}(l) - \textsc{Cost}(l))^{1/t}}{Z} | x, l \in X^+\})$ ▷ Sample a segment to add.
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ ▷ Sample where to move from.
13: $\quad\quad \_, l^- \sim \textsc{Samp.}(\{\langle x, l \rangle : \frac{\sigma(\textsc{Cost}(l) - \textsc{Util.}(l))^{1/t}}{Z} | x, l \in X^-, x = x^+ \wedge \textsc{Util.}(l^-) < \textsc{Util.}(l^+)\})$
14: $\quad\quad R[l^+] \leftarrow R[l^+] \cup \{x^+\};\quad\quad R[l^-] \leftarrow R[l^-] \setminus \{x^+\}$ ▷ Commit transaction.
15: $\quad$ **end if**
16: **end while**;$\quad\quad$ **return** $O$



Figure 3: Illustration of two operations from Algorithm 1. The initial state is on the left. Then, a new segment $x_{89}$ is added to the $l_3$ level. Lastly, the segment $x_0$ is promoted from $l_1$ to $l_2$.

Despite the varying absolute scores of metrics under different references, as explored in Section 4.3, mixing of multiple references leads to an almost perfectly linear combination of the endpoint metric performances. The biggest gains in this respect are obtained by BLEURT, chrF and BLEU, while COMET$_{20}$ is almost unaffected. There is no formal guarantee that the mix of score distributions will not lower the overall Kendall's $\tau$. Nevertheless, a positive conclusion is that if there is budget to only translate 25% of segments with high quality, it should be done and it can only improve the overall evaluation reliability.

**Budget Allocation Algorithm.** We provide a heuristic stochastic algorithm to find an assignment of source segments $S$ to be translated by vendors of different costs and qualities within a specific budget. For the current dataset, we set the cost of a segment in R1, R2, R3, and R4 to be 1, 1, 2, and 3, respectively. Their quality (or "fitness" for the

purpose of automatic evaluation) were set to 1, 2, 4, and 3 based on our observations in Section 4.1. Algorithm 1 contains a hyperparameter $\lambda$ that controls whether the budget will be allocated more towards having multiple references per-segment or more towards having fewer but higher-quality references per-segment, and the temperature $t$ than controls the the sampling randomness.

We formalize the problem with a segment cost $\textsc{Cost}(l)$ for a reference on level $l \in L$ and the utility $\textsc{Util.}(l)$. The levels might correspond to translation vendors which have costs and qualities. In our case, $\textsc{Cost} = \{R1: 1, R2: 1, R3: 2, R4: 3\}$ and $\textsc{Util.} = \{R1: 1, R2: 2, R3: 4, R4: 3\}$. Given a set of source sentences $S$, the goal is to assign the segments to different levels R1 . . . R4. The same segment can be assigned to different quality levels at once, leading to multiple references for that segment. The selection should maximize performance of a particular metric on a number of systems but needs to fit under a fixed budget $B$, i.e.

5

$\sum_l |R_l| \cdot \text{COST}(l) \leq B$. In our setup, to preserve fair comparison, each segment needs to have at least one reference. This is because the smaller the testset, the easier it is to achieve higher but spurious correlations. Therefore, $\bigcup_{l \in L} R_l = S$. The formalization explicitly allows for parts of the testset to be translated multiple times but requires the budget to cover at least the full test set with the cheapest references. This requirement can be fulfilled by subsampling the testset, as commonly done in WMT evaluation campaigns (Kocmi et al., 2023, inter alia).

The pseudocode is provided in Algorithm 1 and explanatory illustration of the two operations in Figure 3. The algorithm continually applies one of the two operations until they can either no longer be applied or the budget is reached. The algorithm will always terminate because ADD increases the cost and utility and PROMOTE increases the utility by at least $\min_{l \in L} \text{COST}(l)$ and $\min_{l_1, l_2 \in L} |\text{UTIL.}(l_2) - \text{UTIL.}(l_1)|$, respectively. Therefore either the budget will be filled or every segment will receive a reference from all vendors.

In Figure 4, we show chrF and COMET$_{20}$ correlations when using the references selected by our algorithm. The optimal preference between quality and quantity changes with increasing budget. Using all of the budget on either quality or quantity would correspond to the bottom or top row, which are not optimal. The best reference configurations for a particular budget, such as $|S| \times 4$, four-times the price of the cheapest translation, contain a mixture of references from R1, R2, R3, and R4 with multiple references for some segments. In addition to the metric correlations in Figure 4, we show the average number of references per source segment in Figure 5. With focus on quality, each segment has fewer references.

## 5. Qualitative Analysis

In Tables 7 and 8, we show a single source segment, one system translation and multiple references and the metric scores. BLEU ranges from 0 to 100 and both extremes are almost achieved just with a different reference. The best human translation led to the lowest BLEU score because of a translation shift. This is not surprising because BLEU operates on the surface-level. Unexpectedly, a similar thing happens also with COMET$_{20}$, which uses a distributed semantic representation of the segments. This shows that parametric model-based metrics are not robust to changes in references. In Table 8, the COMET$_{20}$ difference between references is large due to some translators deciding to drop the verb *"spolupracovat" (collaborate)*, which changes the meaning and the system translation is penalized.



Figure 4: Heatmaps of chrF (left) and COMET$_{20}$ (right) Kendall's $\tau$ correlations on reference configurations created with a specific budget (x-axis) and quality-quantity trade-off $\lambda$ (y-axis). ★ marks the best value in each column (fixed budget). The first column corresponds to the cheapest translation for all test segments, with no room for selection. $\lambda \in [0, 0.7]$ and $t = 0.5$. ◉ With a limited budget, e.g. $2|S|$ or $3|S|$, it makes more sense to add *some* references of a higher quality rather than covering the whole test set with a second reference. With more budget available, multiple references per segment become more beneficial.

## 6. Conclusion

We showed that the quality of references is important for accurate automatic machine translation metrics. The relationship is not straightforward: translatologists' translations, despite being the peak translation quality, are not the best references. Rather, it is the *standard commercial professional translations* that work best for current metrics. The trend applies to both string-matching metrics as well as to parametric model-based ones. Taking the *average over multiple references provides the biggest benefit*, with diminishing returns after 7 references. We also provided a heuristic-based *algorithm for finding a good configuration of references given a budget*, which surpasses optimizing solely for quantity or quantity.

**Future work.** The dataset size prevents system-level investigations. Because there is little point in evaluating segments that are easy to translate, a follow-up approach could prioritize difficult-to-translate segments. This is used by Isabelle et al. (2017) for creating a challenge set. Future works should *quantify* the references quality and ask how many segments are needed to fulfill a certain desideratum, such as effect size or metric accuracy.

**Limitations.** We note the limitation of using a small dataset and a single language translation direction due to the costs of creating multiple rounds of high-quality references. We are convinced the results hold in other scenarios as the effect directions are the same across multiple metrics and setups.

Number of references per segment

Figure 5: Average number of references per one segment allocated by Algorithm 1 with $\tau = 0.5$ (top) and $\tau = 10^{-3}$ (bottom).

| | Metric | R1$^{\text{PE}}$-R1 | R2$^{\text{PE}}$-R2 | R3$^{\text{PE}}$-R3 | R4$^{\text{PE}}$-R4 |
|---|---|---|---|---|---|
| **Layman PE** | BLEU | +0.019 | +0.007 | +0.009 | +0.010 |
| | chrF | +0.027 | +0.015 | +0.016 | +0.019 |
| | TER | +0.026 | +0.015 | +0.013 | +0.014 |
| | COMET$^{20}$ | +0.016 | +0.011 | +0.010 | +0.009 |
| | COMET$^{22}$ | +0.008 | +0.006 | +0.007 | +0.005 |
| | BLEURT | +0.017 | +0.015 | +0.011 | +0.010 |
| **Student PE** | BLEU | +0.010 | +0.001 | −0.004 | −0.001 |
| | chrF | +0.011 | +0.001 | +0.000 | −0.004 |
| | TER | +0.003 | +0.001 | +0.005 | −0.002 |
| | COMET$^{20}$ | +0.010 | +0.003 | −0.002 | −0.002 |
| | COMET$^{22}$ | +0.002 | +0.001 | +0.000 | −0.002 |
| | BLEURT | +0.021 | +0.022 | +0.005 | −0.004 |
| **Prof. PE** | BLEU | +0.035 | +0.011 | +0.007 | −0.000 |
| | chrF | +0.040 | +0.010 | +0.008 | +0.004 |
| | TER | +0.023 | +0.014 | +0.003 | −0.002 |
| | COMET$^{20}$ | +0.016 | +0.006 | +0.000 | +0.005 |
| | COMET$^{22}$ | +0.008 | +0.002 | +0.003 | +0.005 |
| | BLEURT | +0.027 | +0.022 | +0.005 | −0.001 |

Table 6: Difference between using original translations (in Table 2) and post-edited translations as references. Sections are divided based on who did the post-editing (layman, translatology student, or professional translator). This table expands on Table 3. Absolute scores of individual reference subsets are in Table 10.

Table 7: BLEU and COMET$_{20}$ scores of the source "*Three Scottish students named among Europe's best*" and the system translation "*Tři skotští studenti byli zařazeni mezi nejlepší v Evropě*". Both metrics are multiplied by 100. 👁 All references are good translations but the scores vary.

| BLEU | COMET$_{20}$ | Reference |
|---|---|---|
| 10 | 78 | K evropské špičce nově patří i tři skotští studenti |
| 23 | 120 | Tři skotští studenti se umístili mezi nejlepšími v Evropě |
| 23 | 121 | Tři skotští studenti mezi nejlepšími v Evropě |
| 28 | 116 | Tři skotští studenti byli oceněni jako jedni z nejlepších v Evropě |
| 28 | 115 | Tři skotští studenti byli jmenováni jako jedni z nejlepších v Evropě |
| 28 | 114 | Tři skotští studenti byli vyhlášeni jako jedni z nejlepších v Evropě |
| 32 | 117 | Tři skotští studenti byli jmenováni jednimi z nejlepších v Evropě |
| 37 | 122 | Tři skotští studenti byli jmenováni mezi nejlepšími v Evropě |
| 43 | 125 | Tři skotští studenti se zařadili mezi nejlepší v Evropě |
| 43 | 121 | Tři skotští studenti patří mezi nejlepší v Evropě. |
| 43 | 122 | Tři skotští studenti patří mezi nejlepší v Evropě |
| 60 | 127 | Tři skotští studenti zařazeni mezi nejlepší v Evropě |
| 100 | 131 | Tři skotští studenti byli zařazeni mezi nejlepší v Evropě |

Table 8: BLEU and COMET$_{20}$ scores of the source "*Sony, Disney Back To Work On Third Spider-Man Film*" and the system translation "*Disney se vrací, bude spolupracovat se Sony na třetím sólovém Spider-Man filmu*". Both metrics are multiplied by 100. 👁 Some references omit part of the information and COMET$_{20}$ thus penalizes the system translation.

| BLEU | COMET$_{20}$ | Reference |
|---|---|---|
| 4 | -42 | Sony a Disney točí třetí film o Spidermanovi |
| 4 | -33 | Sony a Disney točí třetí film o Spider-Manovi |
| 8 | -9 | Sony a Disney pracují na třetím filmu o Spider-Manovi |
| 8 | -4 | Sony a Disney pokračují v práci na třetím filmu o Spider-Manovi |
| 8 | 1 | Sony a Disney opět pracují na třetím filmu o Spider-Manovi |
| 8 | 15 | Sony a Disney spolupracují na třetím filmu o Spider-Manovi |
| 4 | 16 | Sony a Disney budou spolupracovat při natáčení třetího filmo o Spider-manovi |
| 8 | 28 | Sony a Disney opět spolupracují na třetím filmu o Spider-Manovi |
| 8 | 30 | Sony a Disney budou spolupracovat na třetím filmu o Spider-Manovi |
| 8 | 35 | Sony a Disney budou opět spolupracovat na třetím filmu o Spider-Manovi |
| 17 | 52 | Disney bude znovu spolupracovat se společností Sony na třetím filmu Spider-Man |
| 10 | 64 | Disney bude se Sony dál pracovat na třetím filmu se Spider-Manem |
| 50 | 73 | Disney bude spolupracovat se Sony na třetím sólovém filmu o Spider-Manovi |
| 75 | 99 | Disney se vrací, bude spolupracovat se Sony na třetím filmu o Spider-Manovi |
| 78 | 106 | Disney se vrací, bude spolupracovat se Sony na třetím sólovém filu Spider-Man. |
| 79 | 108 | Disney se vrací, bude spolupracovat se Sony na třetím Spider-Man filmu |
| 100 | 121 | Disney se vrací, bude spolupracovat se Sony na třetím sólovém Spider-Man filmu |

7

|  |  | R1 | R2 | R3 | R4 | R{1,2} | R{1,3} | R{1,4} | R{2,3} | R{2,4} | R{3,4} | R{1,2,3} | R{1,2,4} | R{1,3,4} | R{2,3,4} | Rx |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Average** | BLEU | 0.082 | 0.103 | 0.109 | 0.103 | 0.109 | 0.122 | 0.114 | 0.132 | 0.124 | 0.114 | 0.136 | 0.124 | 0.125 | 0.130 | 0.134 |
| | chrF | 0.090 | 0.125 | 0.128 | 0.123 | 0.121 | 0.135 | 0.124 | 0.148 | 0.139 | 0.135 | 0.146 | 0.135 | 0.140 | 0.147 | 0.147 |
| | TER | 0.082 | 0.092 | 0.114 | 0.105 | 0.095 | 0.120 | 0.107 | 0.125 | 0.117 | 0.120 | 0.121 | 0.110 | 0.123 | 0.127 | 0.124 |
| | COMET$_{20}$ | 0.172 | 0.176 | 0.185 | 0.181 | 0.181 | 0.189 | 0.185 | 0.191 | 0.183 | 0.188 | 0.190 | 0.185 | 0.190 | 0.189 | 0.189 |
| | COMET$_{22}$ | 0.189 | 0.195 | 0.191 | 0.192 | 0.195 | 0.197 | 0.194 | 0.201 | 0.197 | 0.195 | 0.200 | 0.197 | 0.197 | 0.199 | 0.199 |
| | BLEURT | 0.159 | 0.156 | 0.199 | 0.178 | 0.171 | 0.203 | 0.183 | 0.201 | 0.180 | 0.203 | 0.201 | 0.184 | 0.204 | 0.202 | 0.202 |
| | **Average** | 0.129 | 0.141 | 0.154 | 0.147 | 0.145 | 0.161 | 0.151 | 0.166 | 0.157 | 0.159 | 0.166 | 0.156 | 0.163 | 0.166 | 0.166 |
| **Max** | BLEU | 0.082 | 0.103 | 0.109 | 0.103 | 0.116 | 0.122 | 0.118 | 0.135 | 0.132 | 0.111 | 0.138 | 0.137 | 0.121 | 0.135 | 0.137 |
| | chrF | 0.090 | 0.125 | 0.128 | 0.123 | 0.133 | 0.137 | 0.129 | 0.139 | 0.146 | 0.135 | 0.144 | 0.148 | 0.140 | 0.144 | 0.147 |
| | TER | 0.082 | 0.092 | 0.114 | 0.105 | 0.101 | 0.124 | 0.116 | 0.132 | 0.128 | 0.117 | 0.132 | 0.130 | 0.126 | 0.135 | 0.134 |
| | COMET$_{20}$ | 0.172 | 0.176 | 0.185 | 0.181 | 0.177 | 0.184 | 0.185 | 0.180 | 0.183 | 0.183 | 0.181 | 0.183 | 0.184 | 0.181 | 0.182 |
| | COMET$_{22}$ | 0.189 | 0.195 | 0.191 | 0.192 | 0.195 | 0.191 | 0.196 | 0.191 | 0.195 | 0.191 | 0.192 | 0.197 | 0.191 | 0.192 | 0.192 |
| | BLEURT | 0.159 | 0.156 | 0.199 | 0.178 | 0.180 | 0.199 | 0.190 | 0.188 | 0.181 | 0.193 | 0.197 | 0.192 | 0.200 | 0.188 | 0.197 |
| | **Average** | 0.129 | 0.141 | 0.154 | 0.147 | 0.150 | 0.159 | 0.156 | 0.161 | 0.161 | 0.155 | 0.164 | 0.164 | 0.160 | 0.162 | 0.165 |

Table 9: Comparison using either a single or multiple references and taking the average or maximum on segment-level. This table expands on Table 4. The black boxes indicate the reported value of Kendall's $\tau$ visually and are comparable across columns as well as rows.

|  |  | R1[PE] | R2[PE] | R3[PE] | R4[PE] | Rx[PE] | R{1,1[PE]} | R{2,2[PE]} | R{3,3[PE]} | R{4,4[PE]} | R{x,x[PE]} |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Layman PE** | BLEU | 0.101 | 0.111 | 0.117 | 0.113 | 0.144 | 0.092 | 0.107 | 0.113 | 0.108 | 0.140 |
| | chrF | 0.118 | 0.139 | 0.144 | 0.142 | 0.164 | 0.106 | 0.134 | 0.137 | 0.135 | 0.159 |
| | TER | 0.107 | 0.107 | 0.127 | 0.119 | 0.142 | 0.099 | 0.102 | 0.123 | 0.116 | 0.139 |
| | COMET$_{20}$ | 0.188 | 0.187 | 0.195 | 0.190 | 0.198 | 0.183 | 0.184 | 0.193 | 0.188 | 0.197 |
| | COMET$_{22}$ | 0.197 | 0.201 | 0.198 | 0.197 | 0.203 | 0.195 | 0.200 | 0.196 | 0.196 | 0.202 |
| | BLEURT | 0.176 | 0.170 | 0.210 | 0.188 | 0.209 | 0.169 | 0.165 | 0.206 | 0.186 | 0.209 |
| | **Average** | 0.148 | 0.153 | 0.165 | 0.158 | 0.177 | 0.141 | 0.149 | 0.161 | 0.155 | 0.174 |
| **Student PE** | BLEU | 0.092 | 0.104 | 0.105 | 0.102 | 0.123 | 0.089 | 0.107 | 0.108 | 0.103 | 0.130 |
| | chrF | 0.102 | 0.126 | 0.128 | 0.119 | 0.139 | 0.097 | 0.127 | 0.130 | 0.122 | 0.144 |
| | TER | 0.085 | 0.093 | 0.119 | 0.103 | 0.119 | 0.084 | 0.095 | 0.119 | 0.104 | 0.123 |
| | COMET$_{20}$ | 0.182 | 0.179 | 0.183 | 0.179 | 0.186 | 0.179 | 0.179 | 0.186 | 0.181 | 0.188 |
| | COMET$_{22}$ | 0.191 | 0.196 | 0.191 | 0.189 | 0.195 | 0.191 | 0.197 | 0.193 | 0.191 | 0.197 |
| | BLEURT | 0.180 | 0.178 | 0.203 | 0.174 | 0.199 | 0.172 | 0.170 | 0.204 | 0.177 | 0.202 |
| | **Average** | 0.139 | 0.146 | 0.155 | 0.144 | 0.160 | 0.136 | 0.146 | 0.157 | 0.146 | 0.164 |
| **Professional PE** | BLEU | 0.118 | 0.114 | 0.115 | 0.103 | 0.127 | 0.103 | 0.113 | 0.113 | 0.104 | 0.133 |
| | chrF | 0.131 | 0.135 | 0.136 | 0.127 | 0.146 | 0.113 | 0.133 | 0.135 | 0.126 | 0.149 |
| | TER | 0.105 | 0.106 | 0.116 | 0.103 | 0.118 | 0.094 | 0.102 | 0.118 | 0.104 | 0.122 |
| | COMET$_{20}$ | 0.188 | 0.182 | 0.185 | 0.186 | 0.190 | 0.183 | 0.181 | 0.189 | 0.185 | 0.191 |
| | COMET$_{22}$ | 0.198 | 0.198 | 0.195 | 0.196 | 0.199 | 0.195 | 0.199 | 0.195 | 0.195 | 0.200 |
| | BLEURT | 0.186 | 0.178 | 0.204 | 0.176 | 0.197 | 0.177 | 0.172 | 0.206 | 0.179 | 0.202 |
| | **Average** | 0.154 | 0.152 | 0.159 | 0.149 | 0.163 | 0.144 | 0.150 | 0.159 | 0.149 | 0.166 |

Table 10: Metric performance when using post-edited references also jointly with their original versions (averaged at the segment-level). This table expands on Tables 3 and 6.

## Bibliographical References

Loïc Barrault, 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 conference on machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, 1–55. Association for Computational Linguistics.

Ondřej Bojar, Kamil Kos, and David Mareček. 2010. Tackling sparse data issue in machine translation evaluation. In *Proceedings of the ACL 2010 Conference Short Papers*, 86–91, Uppsala, Sweden. Association for Computational Linguistics.

Ondřej Bojar, 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 *Proceedings of the Eighth Workshop on Statistical Machine Translation*, 1–44, Sofia, Bulgaria. Association for Computational Linguistics.

Sheila Castilho, Stephen Doherty, Federico Gaspari, and Joss Moorkens. 2018. Approaches to human and machine translation quality assessment. *Translation quality assessment: From principles to practice*.

Joke Daems and Lieve Macken. 2020. Post-editing human translations and revising machine translations: Impact on efficiency and quality. In *Translation Revision and Post-editing*, 50–70. Routledge.

Andrew M Finch, Yasuhiro Akiba, and Eiichiro Sumita. 2004. How does automatic machine translation evaluation correlate with human scoring as the number of reference translations increases? In *LREC*.

Marina Fomicheva. 2017. *The role of human reference translation in machine translation evaluation*. Ph.D. thesis, Universitat Pompeu Fabra.

Marina Fomicheva, Lucia Specia, and Francisco Guzmán. 2020. Multi-hypothesis machine translation evaluation. In *Proceedings of the 58th Annual Meeting of ACL*.

Markus Freitag, 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.

Markus Freitag, David Grangier, and Isaac Caswell. 2020. BLEU might be guilty but references are not innocent. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 61–71. Association for Computational Linguistics.

Markus Freitag, 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 *Proceedings of the Eighth Conference on Machine Translation*, 578–628. Association for Computational Linguistics.

Markus Freitag, 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 *Proceedings of the Seventh Conference on Machine Translation*, 46–68. Association for Computational Linguistics.

Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021b. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, 733–774. Association for Computational Linguistics.

Jesús González-Rubio, Daniel Ortiz-Martínez, and Francisco Casacuberta. 2012. Active learning for interactive machine translation. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, 245–254. Association for Computational Linguistics.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2016. Can machine translation systems be evaluated by the crowd alone. *Natural Language Engineering*, FirstView:1–28.

Gholamreza Haffari, Maxim Roy, and Anoop Sarkar. 2009. Active learning for statistical phrase-based machine translation. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, 415–423. Association for Computational Linguistics.

Olivier Hamon and Djamel Mostefa. 2008. The impact of reference quality on automatic MT evaluation. In *Coling 2008: Companion volume: Posters*, 39–42.

Pierre Isabelle, Colin Cherry, and George Foster. 2017. A challenge set approach to evaluating machine translation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2486–2496. Association for Computational Linguistics.

Huda Khayrallah, Brian Thompson, Matt Post, and Philipp Koehn. 2020. Simulated multiple reference training improves low-resource machine translation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 82–89. Association for Computational Linguistics.

Věra Kloudová, David Mraček, Ondřej Bojar, and Martin Popel. 2023. Možnosti a meze tvorby tzv. optimálních referenčních překladů: po stopách „překladatelštiny" v profesionálních překladech zpravodajských textů. *Slovo a slovesnost*, 84(2):122–156.

Věra Kloudová, Ondřej Bojar, and Martin Popel. 2021. Detecting post-edited references and their effect on human evaluation. In *Proceedings of the Workshop on Human Evaluation of NLP Systems (HumEval)*, 114–119, Stroudsburg, USA. Association for Computational Linguistics.

Tom Kocmi, 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 *Proceedings of the Eighth Conference on Machine Translation*, 1–42, Singapore. Association for Computational Linguistics.

Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes. 2021. To ship or not to ship: An extensive evaluation of automatic metrics for machine translation. In *Proceedings of the Sixth Conference on Machine Translation*, 478–494. Association for Computational Linguistics.

Tom Kocmi, Vilém Zouhar, Christian Federmann, and Matt Post. 2024. Navigating the metrics maze: Reconciling score magnitudes and accuracies.

Mateusz Krubiński, Erfan Ghadery, Marie-Francine Moens, and Pavel Pecina. 2021. Just ask! evaluating machine translation by asking and answering questions. In *Proceedings of the Sixth Conference on Machine Translation*, 495–506, Online. Association for Computational Linguistics.

Arle Lommel. 2016. Blues for BLEU: Reconsidering the validity of reference-based MT evaluation. *Translation Evaluation: From Fragmented Tools and Data Sets to an Integrated Ecosystem*, 63.

Matouš Macháček and Ondřej Bojar. 2014. Results of the WMT14 metrics shared task. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, 293–301, Baltimore, MD, USA. Association for Computational Linguistics.

Nitin Madnani, Philip Resnik, Bonnie J. Dorr, and Richard Schwartz. 2008. Are multiple reference translations necessary? investigating the value of paraphrased reference translations in parameter optimization. In *Proceedings of the 8th Conference of the Association for Machine Translation in the Americas*, 143–152.

Nitika Mathur, Johnny Tian-Zheng Wei, Markus Freitag, Qingsong Ma, and Ondřej Bojar. 2020. Results of the WMT20 metrics shared task. In *Fifth Conference on Machine Translation - Proceedings of the Conference*, 688–725, Stroudsburg, PA, USA. Association for Computational Linguistics, Association for Computational Linguistics.

Vânia Mendonça, Ricardo Rei, Luísa Coheur, and Alberto Sardinha. 2023. Onception: Active Learning with Expert Advice for Real World Machine Translation. *Computational Linguistics*, 49(2):325–372.

Chenggang Mi, Lei Xie, and Yanning Zhang. 2020. Improving adversarial neural machine translation for morphologically rich language. *IEEE Transactions on Emerging Topics in Computational Intelligence*, (4):417–426.

Myle Ott, Michael Auli, David Grangier, and Marc'Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, 3956–3965. PMLR.

Kishore Papineni, 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 ACL*, 311–318. Association for Computational Linguistics.

Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, 392–395. Association for Computational Linguistics.

Ying Qin and Lucia Specia. 2015. Truly exploring multiple references for machine translation evaluation. In *Proceedings of the 18th Annual Conference of the European Association for Machine Translation*.

Ricardo Rei, 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*, 578–585. Association for Computational Linguistics.

Ricardo Rei, 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 *Proceedings of the Sixth Conference on Machine Translation*, 1030–1040. Association for Computational Linguistics.

Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2685–2702. Association for Computational Linguistics.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of ACL*, 7881–7892. Association for Computational Linguistics.

Zhouxing Shi and Minlie Huang. 2020. Robustness to modification with shared words in paraphrase identification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, 164–171. Association for Computational Linguistics.

Matthew Snover, 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*, 223–231. Association for Machine Translation in the Americas.

Lucia Specia and Kashif Shah. 2018. Machine translation quality estimation: Applications and future perspectives. *Translation quality assessment: from principles to practice*, 201–235.

Marlies van der Wees, Arianna Bisazza, and Christof Monz. 2017. Dynamic data selection for neural machine translation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 1400–1410. Association for Computational Linguistics.

Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 4566–4575.

Giorgos Vernikos, Brian Thompson, Prashant Mathur, and Marcello Federico. 2022. Embarrassingly easy document-level MT metrics: How to convert any pretrained metric into a document-level metric. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, 118–128. Association for Computational Linguistics.

Ying Zhang and Stephan Vogel. 2004. Measuring confidence intervals for the machine translation evaluation metrics. In *Proceedings of the 10th Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages*.

Ying Zhang and Stephan Vogel. 2010. Significance tests of automatic machine translation evaluation metrics. *Machine Translation*, 24:51–65.

Renjie Zheng, Mingbo Ma, and Liang Huang. 2018. Multi-reference training with pseudo-references for neural translation and text generation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 3188–3197. Association for Computational Linguistics.

Vilém Zouhar, Martin Popel, Ondřej Bojar, and Aleš Tamchyna. 2021. Neural machine translation quality and post-editing performance. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 10204–10214. Association for Computational Linguistics.

Vilém Zouhar, Věra Kloudová, Martin Popel, and Ondřej Bojar. 2023. Evaluating optimal reference translations.

# Exploratory Study on the Impact of English Bias of Generative Large Language Models in Dutch and French

**Ayla Rigouts Terryn[1], Miryam de Lhoneux[2]**
[1]KU Leuven, Centre for Computational Linguistics (CCL)
[2]KU Leuven, Department of Computer Science
ayla.rigoutsterryn@kuleuven.be, miryam.delhoneux@kuleuven.be

## Abstract

The most widely used LLMs like GPT4 and Llama 2 are trained on large amounts of data, mostly in English but are still able to deal with non-English languages. This English bias leads to lower performance in other languages, especially low-resource ones. This paper studies the linguistic quality of LLMs in two non-English high-resource languages: Dutch and French, with a focus on the influence of English. We first construct a comparable corpus of text generated by humans versus LLMs (GPT-4, Zephyr, and GEITje) in the news domain. We proceed to annotate linguistic issues in the LLM-generated texts, obtaining high inter-annotator agreement, and analyse these annotated issues. We find a substantial influence of English for all models under all conditions: on average, 16% of all annotations of linguistic errors or peculiarities had a clear link to English. Fine-tuning a LLM to a target language (GEITje is fine-tuned on Dutch) reduces the number of linguistic issues and probably also the influence of English. We further find that using a more elaborate prompt leads to linguistically better results than a concise prompt. Finally, increasing the temperature for one of the models leads to lower linguistic quality but does not alter the influence of English.

**Keywords:** LLM, bias, cross-lingual

## 1. Introduction

In recent years, (generative, pre-trained) large language models (LLMs) have substantially advanced and changed the field of natural language processing (NLP), with large models displaying an "unusually large set of capabilities" (Tamkin et al., 2021) across a wide range of tasks, including acting as a chatbot. Their capabilities and ease of use have contributed to a quick rise in popularity, including among non-expert users. For instance, a recent report on the use of digital technologies in Flanders in 2023 (De Marez et al., 2024) showed that 18% of people in this region use a tool to generate text, music, images, or speech at least monthly. For chatbots specifically, this number drops a little to 14%. Given how recently AI chatbots have become available, this illustrates how fast they are gaining influence.

The undeniably impressive capabilities of the LLMs behind AI chatbots do not imply the technology is without its flaws. For instance, the production of false content by LLMs is common enough that it quickly got a dedicated term: *hallucinations* (see, e.g., Ye et al. (2023)). The models are also known to be *biased* (see, e.g., Vig et al. (2020)). A third issue, which constitutes the central theme of this study, pertains to the *English bias*. This refers to the tendency for LLMs to be predominantly trained on English datasets. The problem goes beyond LLMs, and affects NLP in general: "[e]xisting estimates of how much of top venue NLP research is devoted to English vary a bit, but typically lie in the range of 50-90%" (Søgaard, 2022, p.5254).

The English bias has many effects. Logically, the performance of NLP tools is often best in English. This is clearly illustrated for machine translation, where performance tends to be highest for language pairs that include English, for translation into English, and for English in combination with a closely related language, as illustrated by, e.g., the results of WMT23 (Kocmi et al., 2023). However, this English bias goes beyond performance issues. For instance, De Bruyne (2023) argues that the predominance of English has a (negative) impact on the conceptualisation of emotion detection, as emotions and the ways people verbalise emotions are not universal.

An effect that has not been researched extensively is the linguistic quality of texts generated by LLMs in languages other than English and, specifically, whether and how English bias influences these texts (e.g., presence of anglicisms). The latter is a well-known issue among attentive non-English users of the technology, but very little research can be found where the issue is officially established and analysed. The main goal of this exploratory study is to document general linguistic issues in texts written by generative LLMs and to analyse how often these issues might be traced back to the English bias. The secondary goal is to provide a starting point for future (more extensive) research by testing a methodology based on human annotations and starting to identify the role of some of the main variables, such as the models (and their training data and sizes), languages,

temperature, and prompts.

A brief overview of related research can be found in Section 2. The methodology is described in Section 3, with separate subsections on corpus creation and annotation. Section 4 is dedicated to the findings. Limitations, conclusions, and opportunities for future research are discussed in Section 5.

## 2.    Related Research

The most widely used LLMs like GPT4 (OpenAI et al., 2023) and Llama 2 (Touvron et al., 2023) are trained on large amounts of mostly English data, but are still able to deal with non-English languages (Shi et al., 2023). This English bias leads to lower performance in other languages, especially for low-resource languages (e.g. Hendy et al., 2023, among others) and for tasks that are not translatable (Zhang et al., 2023). This has led researchers to speculate that these models use English as a pivot language in which they reason, prior to generating output in non-English languages. Wendler et al. (2024) empirically test this speculation by inspecting internal model representations via mechanistic interpretability. They develop tasks (translation, repetition and a cloze task) where the output is expected to be in non-English languages (here mainly Chinese, but with controlled experiments in French and Russian) and investigate the latent representations of Llama models at the different layers. They find evidence that the representations in the intermediate layers of these models are closer to English than to other languages, confirming that English may act as a pivot language. Contemporaneously, Zhao et al. (2024) probe LLMs for language-specific information leading them to very similar findings.

Our work focusses on the analysis of the model *outputs*. While Wendler et al. (2024) find that the influence of English in the representation declines to a very small percentage in the last layers of the models, we find clear traces of it in the output. This complements the evidence that English is used as a pivot language in these models. We further contribute a characterisation of *how* English manifests itself in model outputs in non-English languages, here Dutch and French.

## 3.    Methodology

### 3.1.    Corpus

With the intended goal of analysing linguistic output of generative LLMs in non-English languages and the impact of English bias, we decided on a corpus-based approach with expert annotations as the best way to obtain nuanced and fine-grained information. Ideally, a comparable corpus of human-written and LLM-generated texts would allow for a controlled comparison. Because data-driven systems like LLMs perform best on content that is well represented in the training data, we want to work with common text types in the general domain, to avoid adding domain-specific difficulties. At this stage, well-resourced languages that are closely related to English (so, probably some of the *easiest* languages besides English for these models to handle) were preferred. First, because the goal is to include multiple models, we are limited by the availability of models that include the languages. Second, because the quality of texts written by LLMs in low-resource languages can be too low to allow a detailed analyses. Third, because it is worth seeing which problems remain even in such ideal settings. Ideally, lessons learnt from this exploratory study can be used to launch similar studies for low(er) resourced languages.

The languages of choice based on these criteria were Dutch and French, specifically the dialect clusters from Belgium: Flemish and Belgian French. Both are well established and standardised national languages with very limited variation from Dutch and French in the Netherlands and France. An advantage is that it allows a selection of equivalent sources in the same country for both languages. The text collection was based on insights from Schepens et al. (2023) and Muñoz-Ortiz et al. (2023). The former create a German corpus of LLM-generated texts based on prompts including children's books titles, a request to write in German, and for children within a certain age range. The latter use (English) newspaper headlines and the first three words of the article as prompts. Newspaper headlines fit the criteria of general text types that should be well represented in the training data of most models. To create a comparable corpus of publicly available texts in Dutch and French, the online newspapers of the Belgian public broadcasting companies were selected as sources, with articles from the Flemish (Dutch) *Vlaamse Radio- en Televisieomroeporganisatie (VRT)* and from its (Belgian) French counterpart, *Radio-Télévision belge de la Communauté française (RTBF)*.

In terms of the choice of models and settings, there are many potentially influential variables (model architecture, model size, training data size, language distribution in training data, temperature, top P, prompt, etc.). In a small-scale study like this, it is impossible to control for everything and to obtain enough useful data with all potentially relevant combinations of settings. The practical limitations (time and budget) allowed about 500 articles to be annotated. With 50 articles per setting as a reasonable estimate for the minimum volume of text

| | Prompt A | Prompt B |
|---|---|---|
| **Dutch** | Je bent een ervaren journalist bij VRT NWS, de nieuwssite van de Vlaamse openbare omroep. Je moedertaal is Nederlands (Vlaams). Schrijf een artikel voor VRT NWS op basis van volgende titel: [title] | Schrijf een artikel op basis van volgende titel: [title] |
| **French** | Tu travailles en tant que journaliste pour la RTBF, la référence francophone de l'actualité publique belge, et tu as beaucoup d'expérience. Ta langue maternelle est le français (de Belgique). Ecris un article pour la RTBF ayant le titre suivant : [title] | Ecris un article ayant le titre suivant: [title] |
| **English equivalent** | You are an experienced journalist working for [name of broadcasting company], the news website of the [Flemish or Belgian French] public broadcaster. Your native language is [Dutch (Flemish) or French (from Belgium)]. Write an article for [name of broadcaster] based on the following title: [title] | Write an article based on the following title: [title] |

Table 1: Elaborate (A) and concise (B) prompts used in Dutch and French, incl. English translation

required for a meaningful analysis, this amounted to 10 different experimental settings. 50 articles were collected in Dutch and French respectively, spread over various categories of news (national, international, sports, politics, etc.) and making sure the subjects were equivalent in both languages. An overview of the original articles and sources has been added in the appendix. With few exceptions (to find equivalent articles in Dutch and French), only recently published articles were selected to limit the chances of them being included in the training data of the models.

Though we cannot control for all differences between available pretrained models, in the context of this project we looked for (1) one of the largest, best performing models as a reflection of what is currently possible, (2) one (smaller) open source model that allows further research, and (3) one model with more fine-tuning on the non-English language to see whether and how much this can improve results. As *prompt engineering* has also been shown to be influential (White et al., 2023), two different prompts were chosen as additional variables: one elaborate prompt that considers common insights from prompt engineering, like assigning a role (prompt A), and one very concise prompt (prompt B). The exact prompts and an English translation can be found in Table 1. However, as this doubled the number of experiments, to limit the number of articles to 500, the decision was made to only include a fine-tuned (language-specific) model for Dutch, as the lesser-resourced of the two languages. This means the project includes 3 models, all of which are used for Dutch, and two of which are used for French:

- **GPT-4** (OpenAI et al., 2023):

- Settings: used in OpenAI Playground (chat), temperature=1.0, maximum_length=8000, top_P=1.
- Motivation: one of the most powerful and influential models available at the time of the experiment (Zhao et al., 2023).
- Limitations: not open source.

- **Zephyr 7B Beta** (Tunstall et al., 2023):
- Settings: used in the Hugging-Face chat version[1], temperature=0.7, max_new_tokens=1024 (+ click *continue generating* when option is provided after incomplete response), top_P=0.95.
- Motivation: One of the best-performing open source models for Dutch based on (Vanroy, 2023), without specific fine-tuning for Dutch (based on Mistral (Jiang et al., 2023)).
- Limitations: trained on synthetic datasets and more likely to generate problematic content according to the technical report, despite high scores on truthfulness tasks (Vanroy, 2023).

- **GEITje Chat V2 7B** (Rijgersberg and Lucassen, 2023) (only for Dutch):
- Settings: used in LM Studio[2], temperature=2.0, n_predict=-1 ("to allow the model to stop on its own"), top_P=0.95.
- Motivation: open source model specifically fine-tuned for Dutch (also based on Mistral).
- Limitations: no preference optimisation and small for a LLM; GEITje-7B-ultra is superior as a chatbot, but was published after experiments had already started.

---

[1] https://huggingface.co/spaces/HuggingFaceH4/zephyr-chat
[2] https://lmstudio.ai/

| source of articles | | | | av. # | | |
| model | tmp | l | p | tok | typ | typ/tok |
|---|---|---|---|---|---|---|
| GEITje | 0.2 | NL | A | 170 | 77 | 0.59 |
| | | | B | 127 | 66 | 0.67 |
| | 0.85 | | A | 136 | 82 | 0.68 |
| GPT-4 | 1.0 | FR | A | 449 | 233 | 0.52 |
| | | | B | 450 | 217 | 0.48 |
| | | NL | A | 394 | 198 | 0.50 |
| | | | B | 440 | 212 | 0.48 |
| Zephyr | 0.7 | FR | A | 560 | 320 | 0.59 |
| | | | B | 594 | 334 | 0.58 |
| | | NL | A | 494 | 276 | 0.59 |
| | | | B | 528 | 311 | 0.59 |
| **VRT (Dutch)** | | | | 494 | 217 | 0.47 |
| **RTBF (French)** | | | | 441 | 202 | 0.50 |

Table 2: Average (av.) number of tokens (tok), types (typ) (lowercased), and type/token ratio per part of the corpus, distinguishing between model, temperature (tmp), language (l), and prompt (p)

For each model, the default (recommended) settings were selected, except for the maximum length, which was set to the maximum allowed value, so the systems were able to write articles of lengths comparable to those of the original articles. All texts were generated between the 2nd and 31st of January 2024. Because the recommended temperature for GEITje is so much lower than for the other models, some experiments were duplicated using the same settings but a higher temperature (0.85, which is between 1.0 (for GPT-4) and 0.7 (Zephyr)). The result is a collection of 550 articles generated by the LLMs, based on the titles of 50 Dutch and 50 French articles written by human journalists. GEITje had to be stopped manually five times because the systems appeared to be stuck endlessly generating the (exact) same paragraphs. The overview, along with token counts, to indicate the size of the corpus can be found in Table 2. A discussion of these numbers and the type/token ratio can be found in Section 4.

## 3.2. Annotation

### 3.2.1. Annotation scheme

As mentioned, the goal of this project is to establish and document linguistic peculiarities (both clear errors and any text that could be seen as problematic from a linguistic perspective), and to analyse how often issues might be traced back to English. Based on preliminary observations by the leading researcher, an annotation scheme was established to divide these observations into nine categories with labels to allow a nuanced analysis:

- English word/phrase
  - not usually used in Dutch/French
  - sometimes used in Dutch/French
  - very commonly used in Dutch/French
- longer piece of English text
  - part of text
  - entire text
- word/phrase does not exist (*)
- grammar mistake (*)
- spelling mistake (*)
- strange/wrong construction (*)
- strangely used word/phrase (*)
- other linguistic remark
- non-linguistic remark

Options marked with (*) all have three labels:
  - clearly from English
  - could be from English
  - no clear link to English

There are 2 additional markers: 'Not sure' and 'Very minor mistake/humans might write the same'. More detailed information, including examples for each category, can be found in the annotation guidelines.[3] The category for non-linguistic remarks was added to allow annotators to mark strange or non-sensical text passages, even when the issue is not linguistic, but they were instructed to keep this for *meta* information (e.g., the language model writing that it is a language model), or very obviously wrong information that feels weird not to mark (e.g., calling penguins mammals). During the annotation, the annotators did not see the source of the articles, so they could not develop a bias, e.g., when realising that some systems consistently write better or worse texts. All annotations were made in Label Studio (Tkachenko et al., 2020-2022).

### 3.2.2. Annotators

Professional translators with experience translating from English were hired to perform the annotations in their native languages because: (1) translators are assumed to know both their source and target languages very well, (2) translators are supposed to be especially attentive to influences from their source language into their target language, and (3) translators have experience revising and (post-)editing (translated) texts, which can be seen as relevant experience for this task. There were two main annotators: one who annotated all French texts, and one who annotated all Dutch texts. All annotators are native speakers of either Flemish Dutch or Belgian French.

---

### 3.2.3. Inter-annotator agreement

Besides the main annotators, two additional annotators (one professional translator, one researcher with a background in translation; both native speakers) were included to calculate inter-annotator agreement (IAA). The main Dutch annotator and the extra annotators all annotated the same 21 Dutch articles based on the first three Dutch titles.

The first problem with calculating IAA is the lack of a minimum or maximum number of possible annotations, excluding many commonly used metrics. The second problem is that the span selection was not very rigid, both because it can be difficult and not many guidelines were defined in this respect, and because annotators were not always careful about including or excluding trailing spaces. This means that automatic calculations offered within Label Studio (e.g., *basic matching function*: [4]) are quite pessimistic, with agreement scores between 45% and 50%. Therefore, part of the IAA calculation was done manually, examining all annotations and matching them if they were clearly about the same item, even when spans did not overlap perfectly (e.g., annotation of *worst-case scenario's*, or only *worst-case* as English words in Dutch, because the word *scenario* is the same in both languages). The result is a list of 187 possible annotations, with for each possible annotation and annotator an indication of whether the instance was annotated, and, if so, which category was used with which label(s).

This analysis shows good agreement on whether to annotate: annotator pairs agree for 73% to 83% of all 187 items. All three agree on 67% of the items. As this does not consider all of the times where none of the annotators mark anything, this is good agreement. One annotator (not the main one) annotates slightly more than the other two (170 versus 147 and 148 annotations respectively). Out of 62 annotations for which at least one annotator disagrees, 21 are marked as *minor* or *not sure*.

Annotator pairs also agree on which category to assign for 65% to 79% of all 187 instances. The confusion matrices show relatively good agreement overall, with a few logical patterns. One of the matrices is shown in Table 3. The others can be found in the appendix. One annotator is stricter than the others, e.g., annotating wrong punctuation. The most ambiguous categories are *strangely/wrongly used phrase* and *strange/wrong construction*. This was expected, since annotators cannot easily consult resources like dictionaries or grammars to check whether their instinct that a word, phrase, or construction is strange or wrong, is more than a personal preference. However, even seemingly unambiguous categories like *nonexistent word* and

*English word* can be ambiguous, for instance when a Dutch text mentions *gefeed*, i.e., the English word *feed* used with a dutch prefix to conform to Dutch grammar rules. These disagreements are also indications of how the guidelines can be improved in the future, e.g., splitting the rather prescriptive sounding *word/phrase exist* into one category for words/phrases that appear made up by the LLM and have never been written by humans (at least not based on texts that can be found online), and one category for words/phrases that may not be part of the official standard language, but are used by human writers as well.

Most categories include the same 3 labels about potential influence from English, so the agreement on these labels can be compared regardless of the categories. Counting the same labels as perfect agreement, and disagreement with only one point difference as 50% agreement, there was 63% to 77% agreement on the label per annotator pair.

**Dutch vs. French main annotators:** We can get an idea about agreement for French versus Dutch annotations based on the Dutch IAA analysis. No unexpected differences were found, except for the most ambiguous category *strangely/wrongly used word/phrase*. The French annotator used this category a lot more than the Dutch annotator: 18.4 times/1000 tokens on average, versus only 3.6 times/1000 tokens on average. For *spelling mistake* there is a difference of 4.8, and for all other categories, the difference is below 1.5. This is observed in all settings and only for the most ambiguous category, which leads us to conclude that the French annotator was quicker to annotate *strangely/wrongly used words/phrases*, and that this does not necessarily reflect a difference in performance of the LLMs in French versus Dutch. More research is required to confirm this and to improve comparisons across languages. Thus, cross-lingual comparisons in the current project are limited.

In conclusion, agreement is high enough to use the annotations for an exploratory analysis of the texts, provided known disagreements and ambiguities are carefully considered.

## 4. Findings

All analyses are based only on the annotations made by the two main annotators (one per language). Since average text length vary per system, the analysis takes this into account and looks at the number of annotations (per category) per 1000 tokens. This works well, except for the Zephyr model in Dutch, especially with the concise prompt (B), because with this setting, Zephyr wrote 36 of the 50 articles completely in English. In those cases, there will only be a single annotation (*entire text in English*). This makes it seem as if there are

---

[4]https://docs.humansignal.com/guide/stats

| annotator A → vs B ↓ | English word/phrase | grammar mistake | longer piece of English text | non-linguistic remark | other linguistic remark | spelling mistake | strange/wrong construction | strangely/wrongly used word/phrase | word/phrase does not exist | #NA | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **English word/phrase** | 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | **15** |
| grammar mistake | 0 | 23 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 11 | **37** |
| **longer piece of English text** | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | **2** |
| **non-linguistic remark** | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | **1** |
| **other linguistic remark** | 0 | 0 | 0 | 4 | 5 | 0 | 0 | 0 | 0 | 4 | **13** |
| **spelling mistake** | 1 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 5 | **18** |
| **strange/wrong construction** | 0 | 0 | 0 | 0 | 0 | 0 | 23 | 0 | 0 | 11 | **34** |
| **strangely/wrongly used word/phrase** | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 26 | 0 | 5 | **33** |
| **word/phrase does not exist** | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 14 | 1 | **17** |
| **#NA** | 1 | 3 | 0 | 0 | 0 | 6 | 1 | 3 | | 3 | **17** |
| **Total** | **15** | **26** | **2** | **5** | **5** | **19** | **28** | **31** | **16** | **40** | **187** |

Table 3: Confusion matrix based on the annotations of two of the annotators

| | GPT-4 temp:1 | | | | GEITje temp:.2 | | GEITje temp:.85 |
|---|---|---|---|---|---|---|---|
| **av. # annotations per category, per 1000 tokens** | **FR** | | **NL** | | **NL** | | **NL** |
| | **A** | **B** | **A** | **B** | **A** | **B** | **A** |
| English word/phrase | 1.23 | 1.22 | 2.39 | 1.68 | 2.82 | 1.34 | 1.74 |
| word/phrase does not exist | 0.19 | 0.26 | 0.51 | 0.49 | 0.18 | 0 | 0.28 |
| grammar mistake | 2.47 | 2.43 | 1.94 | 2.63 | 2.25 | 2.10 | 2.86 |
| spelling mistake | 2.55 | 2.55 | 4.91 | 5.74 | 8.15 | 13.26 | 10.73 |
| strange/wrong construction | 2.66 | 3.02 | 2.21 | 3.10 | 2.01 | 1.75 | 4.36 |
| strangely/wrongly used word/phrase | 14.54 | 15.09 | 2.53 | 2.47 | 0.45 | 0.50 | 1.70 |
| other linguistic remark | 0.45 | 0.27 | 1.02 | 0.98 | 0.37 | 0.45 | 0.71 |
| non-linguistic remark | 0.89 | 0.55 | 0.57 | 0.89 | 2.85 | 1.07 | 4.61 |
| **all annotations (excl. non-ling.)** | 24.08 | 24.83 | 15.50 | 17.10 | 16.22 | 19.40 | 22.37 |
| **all annotations** | 24.97 | 25.38 | 16.08 | 17.98 | 19.07 | 20.47 | 26.98 |
| **text written completely in English** | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| **average % of annotations with:** | | | | | | | |
| clear English influence | 7% | 6% | 8% | 13% | 24% | 6% | 4% |
| potential influence from English | 36% | 39% | 14% | 26% | 26% | 24% | 33% |
| no clear influence from English | 57% | 54% | 78% | 60% | 50% | 70% | 63% |

Table 4: Averaged findings per setting (language FR or NL; prompt A or B) of GPT-4 and GEITje (with recommended temperature of .2, then with temperature of .85)

very few annotations in the other categories in this setting (because these cannot be annotated in the English texts), which is not representative (the few Dutch texts do contain a lot of annotations). Therefore, this setting is often excluded from the general analyses.

**Number of tokens and types:** A first observation based on the information in Table 2 is that the average lengths of articles differs substantially. GPT-4's average article length is closest to that of the original articles. GEITje regularly writes articles that consist just of (a rephrasing of) the original title (28 of the 150 articles written by GEITje have <50 words). The type/token ratio is also similar for the original articles and the ones written by GPT-4, but higher for Zephyr and GEITje, indicating those models use a more diverse vocabulary. This is especially noteworthy given the fact that annotators indicated that the generated articles were very repetitive. As mentioned, GEITje was even stopped five times because the systems appeared stuck endlessly generating the exact same paragraphs.

**Zephyr:** As expected (because it is a smaller model than GPT-4 and not specifically fine-tuned on Dutch like GEITje), the linguistic quality of texts written by Zephyr is clearly the worst out of the three models. As mentioned, it systematically (36 out of 50 prompts) writes an English article when prompted with the concise prompt in Dutch. It does so for the French concise prompt four times as well, and also twice for the Dutch elaborate prompt. Interestingly, this happens less in French than in Dutch, but with the French prompts, there were also two articles written completely in German and one in Spanish. When writing in the expected language, there are still regularly longer pieces of texts written in English in all settings (20 times in 200 articles). The text written in the expected language contains more annotations on average than the texts written by the other models. For every 1000 tokens, there are on average 40 (French, prompt A), 38 (French, prompt B), and 59 (Dutch, prompt A) annotations, compared to 25 average across the other models. There are especially many *strangely/wrongly used word/phrase* annotations, and, in Dutch, a lot of *word/phrase does not exist* annotations (11 such annotations per 1000 tokens). The proportion of those annotations where an influence of English is expected is not much higher than for the other models: 10% clearly suspected influence and 65% no suspected influence on average in French, and 20% and 66% respectively in Dutch with prompt A. Since Zephyr is much worse than the other two models, the following analyses focus mainly on GPT-4 and GEITje.

**GPT-4 vs. GEITje (Dutch):** Both GPT-4 and GEITje perform a lot better than Zephyr, with fewer annotations on average and fewer texts written in English. The average number of annotations per 1000 tokens (per category) can be seen in Table 4, as well as the proportion of annotations where an influence from English is suspected. When comparing the performance in Dutch of GPT-4 and GEITje (with recommended temperature of .2), a few interesting observations can be made. First, despite GEITje's fine-tuning on Dutch, the experimental setting in Dutch with fewest linguistic annotations was using GPT-4 with Prompt A, though closely followed by GEITje with Prompt A. When non-linguistic remarks are included, GEITje falls further behind. This leads us to a first tentative conclusion regarding this comparison: fine-tuning on Dutch has improved the linguistic quality of GEITje such that it can compete with a much larger LLM like GPT-4 that is not specialised in Dutch. The fact that GEITje is based on the same model (Mistral) as Zephyr, which performs much worse, further strengthens this conclusion. However, the overall non-linguistic quality of texts written by GEITje is not comparable to GPT-4 yet. This is not just reflected in the explicitly annotated *non-linguistic remarks*, but also in the comments shared by the annotators, e.g., about how repetitive the articles written by GEITje are.

**English vs. French (GPT-4):** Another observation is that there are many more annotations in GPT-4's texts written in French than in Dutch, but, as discussed in the previous section, this can largely be attributed to a disagreement between the Dutch and French annotators on how quickly to use the category *strangely/wrongly used word/phrase*. Considering some room for annotator disagreement in the cross-lingual analysis, it is actually remarkable how similar the average number of annotations are per category in both languages. In terms of the suspected influence of English, more research is needed with cross-lingual comparisons, but this influence appears more present in Dutch than in French. In Dutch, there are proportionally slightly more annotations with a clear suspected influence from English, and one text written completely in English instead of Dutch. This is in line with the findings for Zephyr, though with a better average (linguistic) quality.

**Prompt A vs. Prompt B:** Across models, the elaborate prompt (A) leads to linguistically better results than the concise prompt (B), but the difference is not always significant. It is most striking for Zephyr in Dutch, where prompt B leads to 36/50 texts written completely in English, and prompt A prevents this from happening in all but 2 cases. The two times where texts were written completely in English by the other two models was also with the concise prompt. For GPT-4 and GEITje respectively, there are on average 1.59 and 3.81 more linguistic annotations per 1000 tokens for the same experiments with prompt B instead of A in Dutch.

The difference is even smaller for French at 0.75. This influence does not appear to affect any specific type of annotations more than the others, and though the general improvement with the elaborate prompt is consistent, it is not statistically significant according to a paired t-test.

**Temperature:** Because the recommended temperature for GEITje (.2) is much lower than that of Zephyr (.7) and GPT-4 (1.0), GEITje was also tested with a higher temperature. This substantially increased the number of linguistic and non-linguistic annotations. There are significantly (paired t-test, $p < 0.05$) more *strangely/wrongly used word/phrase* annotations with the higher temperature. It is also the setting with most non-linguistic annotations per 1000 tokens out of all, and annotators comment more about strange hallucinations in this setting. The influence of English does not appear affected by the temperature. A notable example of nonsensical output was the following response to prompt A: "Dit is niet mogelijk, aangezien ik een AI-assistent ben die geen Nederlands spreekt." The English translation of this response in Dutch reads: "This is impossible, since I am an AI assistant who does not speak Dutch."

**Influence of English**: Averaged over all settings, 16% of the annotations are labelled as clearly influenced by English. No influence was suspected for 61% of the annotations. There are big differences per setting and category, but since there are sometimes only a few annotations of a category in a setting, the analysis is limited to those where the differences are large and consistent enough to indicate possible generalisation. Curiously, when averaging over all categories, GEITje displays both the highest and lowest percentage of annotations with a clearly suspected influence from English: 30% for prompt A and the recommended low temperature, versus only 4-5% in the other two settings. On closer inspection, the higher number appears due to a few repeated instances that have a big effect because GEITje's texts tend to be short and contain few annotations. For instance, in one text, *Tour of California* is repeated six times and consistently tagged as clearly influenced by English.

Apart from such cases, the overall influence of English in texts generated by GEITje does appear less obvious than with the other models. Analysing this influence per category results in a few more interesting observations.

The *English word/phrase* annotations regularly concern words that are also often used by native speakers of Dutch and French, except for the texts generated by Zephyr, and the French texts generated by GPT-4, where an average of 75% of those annotations are labelled as not generally used in French or Dutch. This is much lower in the other experiments (combined average of 16%). With grammar and spelling mistakes, there is very little suspected influence of English (on average only 3% with a clear reported influence).

A larger percentage is seen for the *word/phrase does not exist* category (see also the section on IAA for a discussion about this category). Zephyr "makes up" a lot of words, with up to 10.7 such annotations per 1000 words in Dutch using prompt A. Some of the annotations in this category consist of seemingly literal "translations" of English words or phrases. For instance, when referring to traffic congestion, the Dutch word *verkeerscongestatie* is used, which does not exist (0 hits when Googling this word). The first part, *verkeer*, is a correct equivalent of *traffic*. The *s* is correctly added for a correct compound. but *congestatie* is an adaptation of *congestion* that may look Dutch, but does not exist as such (the equivalent of *congestion* can be *congestie* in some cases, but not *congestatie*). And even if *congestatie* were the correct term in Dutch, the compound of with *verkeer* does not exist. Instead, the word *file* is used to refer to traffic congestion. Similarly, in French the phrase *si vous ne pouvez pas les battre, alors rejoignez-les* is used (from *if you cannot beat them, join them*). This French phrase has been used online before (10 Google hits), but is a clear anglicism.

**Other observations:** Another noteworthy observation made by the annotators was that the writing was inconsistent. This was true for spelling (e.g., *rechts-extremisme* and *rechtsextremisme* in the same article), vocabulary (e.g., switching between *materieel* and *materiaal* in the same article), and punctuation (e.g., French « and English " quotation marks in the same article). Often, multiple options can be considered correct, but it is good practice to remain consistent within a single text. However, since these models are trained on many different types of texts (the exact training data is not disclosed), and don't necessarily contain information about the boundaries between different texts in the training data, it is not surprising that the output contains some inconsistent writing.

As a concluding remark, it is interesting to see annotators comment on the *stylistic features* of generated texts.

> "Certain stylistic features often demonstrate the intervention of artificial intelligence, such as the logical connectors between parts of articles (*en somme, en conclusion, en conséquence, ...*) which are too obvious, unnatural and which would be more nuanced or subtle in a classic article. What also stands out, for being unnatural, is the emphasis often used to describe a situation, a use of dramatic adjectives to describe a sometimes banal situation in an attempt to add effect, I

guess, but it doesn't work at all."

## 5. Limitations and Conclusions

This exploratory study is a first step towards documenting and better understanding the linguistic qualities of LLMs when writing in Dutch and French, with special focus on the common English bias that is due to the relative overrepresentation of English in the training data o most LLMs. To this end, articles were generated by three different models based on real newspaper headlines, and the resulting corpus was annotated by professional translators for linguistic errors and peculiarities.

Model, language, prompt, and temperature all have a clear impact on results. The difference is noticeable when looking at a simple surface measure like type/token ratio, which is especially high for GEITje, despite repetitive texts. Zephyr is clearly outperformed by the other two models. The most striking result of Zephyr is the number of texts written completely in English instead of Dutch, and the fact that out of the 100 articles to be written by Zephyr based on French prompts French, four were written in English, two in German, and one in Spanish. Linguistically, both GPT-4 and GEITje perform much better and show relatively similar results, indicating that fine-tuning on a specific language can compensate for a smaller model in terms of linguistic quality.

Cross-lingual analyses indicate that the linguistic quality is better in French than Dutch. Comparing a concise and a more elaborate prompt reveals an increased linguistic quality for the latter, though the size of the impact varies per model. Increasing GEITje's very low recommended temperature reduces linguistic quality and increases the number of non-linguistic remarks.

The influence of English is clearly seen for 16% of the annotations on average and can be illustrated very clearly when words or phrases appear to be literally translated from English into Dutch or French words or phrases that are (almost) never used by native speakers.

The main limitations of this study are (1) its scale (limited amount of data per experimental setting), (2) the limited number of languages (only well-resourced languages that are closely related to English), and (3) the potential ambiguity of the annotations. However, the findings can help to narrow down research questions and improve methodologies for experiments on a larger scale. The annotation scheme should be refined to reduce the ambiguity and allow more cross-lingual comparisons.

Since some findings were already relatively clear even with the current setup (e.g., positive impact of elaborate prompt, especially for smaller model),

future research can focus more on, e.g., cross-lingual experiments or fine-grained comparison of annotation categories. Given these improvements, expanding the experiments to include more languages will help to improve our understanding of the linguistic qualities of this influential technology. Another worthwhile direction for future research would be to expand the experiments to include more diverse (and perhaps less formal) text types, as the current setup only covered news articles. Further research could also be dedicated to relating and comparing these findings to human linguistic transfer. Knowing whether the influence of human L1 on L2 is similar to the English bias exhibited by LLMs can help to better understand and predict the performance of LLMs.

## 6. Bibliographical References

Luna De Bruyne. 2023. The Paradox of Multilingual Emotion Detection. In *Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis*, pages 458–466. Association for Computational Linguistics.

Lieven De Marez, R Sevenhant, F Denecker, A Georges, G Wuyts, and D Schuurman. 2024. Imec.digimeter.2023. Digitale trends in Vlaanderen.

Amr Hendy, 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.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, prefix=de las useprefix=false family=Casas, given=Diego, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B.

Tom Kocmi, 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 *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42. Association for Computational Linguistics.

Alberto Muñoz-Ortiz, Carlos Gómez-Rodríguez, and David Vilares. 2023. Contrasting Linguistic Patterns in Human and LLM-Generated Text.

OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo 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, 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, Giambattista 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, Pokorny, Michelle Pokrass, Vitchyr 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 Sheppard, 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 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. 2023. Gpt-4 technical report.

Edwin Rijgersberg and Bob Lucassen. 2023. Geitje: een groot open nederlands taalmodel.

Job Schepens, Nicole Marx, and Benjamin Gagl. 2023. Can we utilize Large Language Models (LLMs) to generate useful linguistic corpora? A case study of the word frequency effect in young German readers.

Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. Language models are multilingual chain-of-

thought reasoners. In *The Eleventh International Conference on Learning Representations*.

Anders Søgaard. 2022. Should We Ban English NLP for a Year? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5254–5260. Association for Computational Linguistics.

Alex Tamkin, Miles Brundage, Jack Clark, and Deep Ganguli. 2021. Understanding the Capabilities, Limitations, and Societal Impact of Large Language Models.

Maxim Tkachenko, Mikhail Malyuk, Andrey Holmanyuk, and Nikolai Liubimov. 2020-2022. Label Studio: Data labeling software. Open source software available from https://github.com/heartexlabs/label-studio.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. Zephyr: Direct distillation of lm alignment.

Bram Vanroy. 2023. Language Resources for Dutch Large Language Modelling.

Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. In *Advances in Neural Information Processing Systems*, volume 33, pages 12388–12401. Curran Associates, Inc.

Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2024. Do Llamas Work in English? On the Latent Language of Multilingual Transformers. ArXiv:2402.10588 [cs].

Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C. Schmidt. 2023. A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT.

Hongbin Ye, Tong Liu, Aijia Zhang, Wei Hua, and Weiqiang Jia. 2023. Cognitive Mirage: A Review of Hallucinations in Large Language Models.

Xiang Zhang, Senyu Li, Bradley Hauer, Ning Shi, and Grzegorz Kondrak. 2023. Don't trust chatGPT when your question is not in english: A study of multilingual abilities and types of LLMs. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A Survey of Large Language Models.

Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji Kawaguchi, and Lidong Bing. 2024. How do large language models handle multilingualism?

## 7. Appendix

### 7.1. Original articles from VRT & RTBF

Tables 5 and 6 list the original articles from the websites of VRT (Flemish) and RTBF (Belgian French) respectively.

| nr | title | author | pub. date |
|---|---|---|---|
| 1 | Vogelgriep treft voor het eerst ijsbeer op Noordpool: "Hier hebben we geen handleiding voor" | Stien Schoofs | 03/01/2024 |
| 2 | Taiwan ontdekt drie Chinese ballonnen in de buurt van luchtmachtbasis | Veerle De Vos | 03/01/2024 |
| 3 | Drie Belgische drugsuithalers opgepakt in Rotterdamse haven, jongste amper 14 jaar | Victor Van Driessche, Belga | 03/01/2024 |
| 4 | Onderzoekers gaan kwab-alen tellen in ondergelopen weides aan Grote Nete | Radio 2, Mathieu Verstichel | 03/01/2024 |
| 5 | Wil je echt vermageren? Zeg dan niet "350 kcal", maar wel "een halfuurtje fietsen" | Dominique Fiers | 02/01/2024 |
| 6 | Rector universiteit Harvard stapt op na ophef over aanpak van antisemitisme en beschuldiging van plagiaat | Nils Schillewaert | 02/01/2024 |
| 7 | Deel van parcours in Gullegem staat onder water: "Maar de veldrit komt niet in het gedrang" | not mentioned | 03/01/2024 |
| 8 | Waarom vond je Belgische tomaten in de winkelrekken op reis in Spanje en Griekenland? | Dennis van den Buijs | 03/01/2024 |
| 9 | Opnieuw miljoenen extra fietsers geteld in provincie Antwerpen: "Alle overheden samen moeten moordstrookjes aanpakken" | Radio 2, Mathieu Verstichel | 03/01/2024 |
| 10 | "Schommelmoment" van verkeersanker Mona krijgt trofee voor mooiste Radio2-moment van 2023 | Radio 2, Martijn Donné | 02/01/2024 |
| 11 | Vliegtuigje neergestort tegen geparkeerde auto in Spa: piloot en inzittende overleden | Belga, Kirsten Sokol | 28/01/2024 |
| 12 | New York Times: "Tijdelijk staakt-het-vuren in Gaza van twee maanden in de maak" | Freek Willems | 28/01/2024 |
| 13 | Oekraïense geheime dienst ontdekt fraude bij wapenaankoop, bijna 37 miljoen euro verdwenen | Freek Willems | 28/01/2024 |
| 14 | Tien landen schorten financiering VN-agentschap UNRWA op na beschuldigingen over betrokkenheid bij terreuraanval Hamas | Kirsten Sokol, Joris Truyts, Freek Willems | 27/01/2024 |
| 15 | Van Taylor Swift over Celine Van Ouytsel tot Emma Watson: "deepnudes" overspoelen het internet (en niet alleen op X) | Maarten Bockstaele | 28/01/2024 |
| 16 | Waarom de landbouwers in Europa en bij ons actievoeren | not mentioned | 28/01/2024 |
| 17 | Frans gerecht verklaart acteur Alain Delon beperkt handelingsonbekwaam | Lina El Bakkali, Belga | 28/01/2024 |
| 18 | Intermittent fasting blijft een hype, werkt het ook? | Radio 1, Maxine Rappé | 28/01/2024 |
| 19 | Koning Charles III maakt het goed na zijn prostaatbehandeling | Lukas Lecluyse | 26/01/2024 |
| 20 | Meer vaders nemen een halve dag per week ouderschapsverlof: "Heeft minder impact op je werkweek en op je loon" | Sandra Cardoen | 27/09/2023 |
| 21 | Sport- en energiedrankjes Prime zijn hype bij jongeren, maar hoe ongezond zijn ze? | Wim De Maeseneer, Nils Schillewaert | 04/08/2023 |
| 22 | Klassieke muziek verbindt ons: zelfs onze hartslag synchroniseert | Radio 1, Maxine Rappé | 10/11/2023 |
| 23 | Minister Tinne Van der Straeten ziet geen reden om snel over nieuwe abortuswet te stemmen: "Thema verdient beter" | Joris Truyts, Nils Schillewaert | 27/01/2024 |
| 24 | Duizenden deelsteps verdwijnen uit Brusselse straatbeeld | BRUZZ, Emmanuel Vanbrussel | 23/01/2024 |
| 25 | Nog drie weken tot oudejaar, maar we weten het nu al zeker: 2023 wordt warmste jaar ooit gemeten | Vincent Merckx | 06/12/2023 |
| 26 | Yana's (21) eetstoornis verergerde door TikTok: bijna helft van jongeren ziet berichten over diëten en mager zijn | Dorien Vanmeldert | 07/10/2023 |

| nr | title | author | pub. date |
|---|---|---|---|
| 27 | Apple stoot Samsung na 12 jaar van de troon als grootste smartphoneverkoper ter wereld | Lukas Lecluyse | 17/01/2024 |
| 28 | Oudste bos ooit van 385 miljoen jaar oud strekte zich uit over 400 kilometer | Michaël Torfs | 13/01/2024 |
| 29 | Opnieuw tienduizenden Duitsers op straat tegen uiterst rechts | Joris Truyts, Belga | 27/01/2024 |
| 30 | Batopin vindt moeilijk locaties voor geldautomaten: "Alle suggesties zijn welkom" | Radio 2, Fred Breuls, Bente Vandekeybus | 30/01/2024 |
| 31 | "Hatsjie": het hooikoortsseizoen is begonnen, ontdek op onze pollenbarometer welke pollen je moet vrezen | Vincent Merckx, Belga | 30/01/2024 |
| 32 | Twee slachtoffers door storm Isha in Verenigd Koninkrijk, tienduizenden huishoudens zonder stroom in Ierland | Ellen Maerevoet, Maarten Bockstaele, Sara Van Poucke, Belga | 22/01/2024 |
| 33 | Tot -48 graden (en het voelt nóg kouder): Vlamingen getuigen over ijzige kou in Canada | Zico Saerens | 13/01/2024 |
| 34 | 22 Genkse basisscholen hebben eigen bibliotheek: "We willen duidelijk maken dat lezen overal kan" | Radio 2, Fred Breuls | 22/12/2023 |
| 35 | CHECK - Ja, een loonsverhoging levert op voor de staatskas, zoals PS-voorzitter Paul Magnette zegt, maar er zijn ook extra kosten | Nele Baeyens, RTBF, Dorien Vanmeldert | 23/01/2024 |
| 36 | 22-jarige Van Uden klopt Groenewegen en Merlier op weg naar eerste sprintzege | not mentioned | 30/01/2024 |
| 37 | Neuralink plaatst eerste hersenimplantaat in menselijk proefpersoon: "We staan nog veraf van hacken van gedachten" | Chris Van den Abeele, Belga, Pieterjan Huyghebaert | 30/01/2024 |
| 38 | Baby "van nog geen uur oud" gevonden in boodschappentas in Londen | Freek Willems | 19/01/2024 |
| 39 | Brand verwoest al bijna 600 hectare van beschermd natuurpark in Argentinië | Lina El Bakkali, Belga | 28/01/2024 |
| 40 | Wilm Vermeir verkozen tot Ruiter van het Jaar, ook zijn paard IQ van het Steentje valt in de prijzen | niet vermeld | 16/01/2024 |
| 41 | Wallonië spendeert per inwoner 70 procent meer aan openbaar vervoer dan Vlaanderen | Rik Arnoudt | 27/01/2024 |
| 42 | Met ChatGPT en geleende spikes: het knotsgekke olympische succesverhaal van John Heymans | Sporza | 29/01/2024 |
| 43 | Mali, Burkina Faso en Niger trekken zich terug uit ECOWAS-verbond | Maarten Bockstaele | 29/01/2024 |
| 44 | Japanse maanlander werkt opnieuw, meer dan een week na de landing | Kathleen Heylen | 29/01/2024 |
| 45 | Eén dode bij aanval van gewapende en gemaskerde mannen in kerk in Istanbul | Joris Truyts | 28/01/2024 |
| 46 | Pakistan voert luchtaanvallen uit op Iran, vrees voor escalatie in de regio | Sara Van Poucke, Nils Schillewaert | 18/01/2024 |
| 47 | Na 2 jaar zicht op nieuwe regering in Noord-Ierland, mét voor het eerst premier van Sinn Féin | Freek Willems | 30/01/2024 |
| 48 | Tomorrowland maakt line-up bekend: op de affiche onder meer David Guetta, Dimitri Vegas & Like Mike en Amber Broos | Belga | 25/01/2024 |
| 49 | Amerikaanse krant The New York Times klaagt OpenAI en Microsoft aan, omdat ze miljoenen artikels gebruikt hebben om ChatGPT te trainen | Wim De Maeseneer, Belga | 27/12/2023 |
| 50 | Drugsdealer loopt tegen de lamp in Brussel, probeert agenten in burger drugs te verkopen | Radio 2, Evi Walschaers | 30/01/2024 |

Table 5: VRT articles

| nr | title | author | pub. date |
|---|---|---|---|
| 1 | Grippe aviaire : un ours polaire infecté en Alaska, une première | Johanne Montay | 08/01/2024 |
| 2 | Taïwan : à quatre jours des présidentielles, le lancement d'un satellite chinois provoque des messages d'alerte | La rédaction, Belga | 09/01/2024 |
| 3 | Rotterdam : arrestation d'un baron de la drogue recherché par la Belgique | Belga, Alain Lechien | 05/01/2023 |
| 4 | Ecraser les oursins violets au marteau pour sauver l'écosystème marin en Californie | Laurick Ayoub sur base d'un reportage de Philippe Jacquemotte | 28/12/2023 |
| 5 | Pourquoi faut-il continuer à faire du sport en hiver ? | Aurélien David via La Une | 20/11/2023 |
| 6 | Suite à plusieurs polémiques, la présidente d'Harvard annonce sa démission | La rédaction | 02/01/2024 |
| 7 | Michael Vanthourenhout s'impose en solitaire à Gullegem en l'absence du "Big Three" | Jâd El Nakadi avec Belga | 06/01/2024 |
| 8 | Selon l'observatoire des prix, 60% des produits alimentaires coûtent moins cher en Belgique qu'ailleurs | QR l'actu | 08/01/2024 |
| 9 | Liège : mauvais bilan 2023 en matière de progrès pour la mobilité cyclable | Marie Bourguignon | 02/01/2024 |
| 10 | Julie Compagnon, les habitants de Bertrix et… la police ont explosé les décibels pour Viva for Life | Par Viva for Life via La Une | 22/12/2023 |
| 11 | Spa : deux morts dans le crash d'un petit avion de tourisme près de l'aérodrome | Olivier Genon | 28/01/2024 |
| 12 | Guerre au Proche-Orient : de violents affrontements sont en cours aux abords des deux principaux hôpitaux de Khan Younès à Gaza | Par La rédaction Info avec Belga | 27/01/2024 |
| 13 | Détournement de 40 millions de dollars par des responsables militaires et chefs d'entreprise ukrainiens | Par La rédaction Info avec Belga | 28/01/2024 |
| 14 | Guerre Israël-Gaza : l'aide à l'Unrwa déjà suspendue par sept pays | Par la rédaction avec AFP | 27/01/2024 |
| 15 | "Protégez Taylor Swift" : les fans se mobilisent pour la défendre contre des deepfakes pornographiques | Par Eléna Lefèbvre | 26/01/2024 |
| 16 | Que compte faire le monde politique en réponse au mécontentement des agriculteurs ? | BELGA – ERIC LALMAND | 28/01/2024 |
| 17 | France : Alain Delon placé sous sauvegarde de justice | Par la rédaction avec AFP | 28/01/2024 |
| 18 | Pour perdre du poids, mieux vaut prendre son petit-déjeuner à 11 heures | Par RTBF avec AFP | 20/06/2022 |
| 19 | Royaume-Uni : le roi Charles III quitte l'hôpital après une opération de la prostate | Par la rédaction avec AFP | 28/01/2024 |
| 20 | Le congé parental n'a jamais été aussi populaire qu'en 2023 en Belgique | Par la rédaction avec Belga | 28/01/2024 |
| 21 | Troubles du sommeil : les boissons énergisantes mises en cause, même à petites doses | Par RTBF avec ETX | 28/01/2024 |
| 22 | La pratique d'un instrument de musique et du chant améliorerait la santé cérébrale des personnes âgées | Par ETX Daily Up édité par Céline Dekock | 30/01/2024 |
| 23 | Avortement : le chantage conservateur du CD&V | Par Philippe Walkowiak | 30/01/2024 |
| 24 | Trottinettes partagées à Bruxelles : Uber et Voi, opérateurs recalés, attaquent la Région en justice | Par Karim Fadoul | 30/01/2024 |
| 25 | Le record de température de 48,8°C en Europe continentale confirmé par l'ONU | Par Marine Lambrecht | 30/01/2024 |
| 26 | Legging legs : la nouvelle tendance controversée et dangereuse qui glorifie la maigreur | Par RTBF avec ETX | 30/01/2024 |
| 27 | Apple dépasse Samsung pour la première fois sur le marché des smartphones | Par Anthony Mirelli | 17/01/2024 |

| nr | title | author | pub. date |
|---|---|---|---|
| 28 | Des scientifiques pensent avoir découvert la plus vieille forêt du monde | Par RTBF Tendance avec AFP | 22/12/2019 |
| 29 | Des milliers de personnes manifestent à nouveau contre l'extrême droite en Allemagne | Par la rédaction avec Belga | 27/01/2024 |
| 30 | La Belgique maintiendra l'accessibilité au cash et aux agences bancaires | Par Maud Wilquin | 25/01/2024 |
| 31 | Les premiers pollens de l'année sont arrivés : la saison des allergies a officiellement commencé | Par Marine Lambrecht | 30/01/2024 |
| 32 | Tempête Isha : un mort en Ecosse, fortes perturbations en Irlande | Par la rédaction avec Belga | 22/01/2024 |
| 33 | Une vague de froid fait au moins 50 morts aux États-Unis | Par la rédaction info avec Belga | 20/01/2024 |
| 34 | 20 minutes de lecture obligatoire, tous les vendredis, au lycée François de Sales à Gilly | Par Simon Gerard | 30/01/2024 |
| 35 | Une augmentation des salaires de 2% permet-elle de réduire le déficit de l'État de deux milliards, comme l'affirme Paul Magnette ? | Par Grégoire Ryckmans avec nws check VRT | 23/01/2024 |
| 36 | Casper van Uden surprend Dylan Groenewegen et Tim Merlier sur la première étape de l'AlUla Tour | Par Cédric Lizin | 30/01/2024 |
| 37 | Elon Musk annonce que Neuralink a posé son premier implant cérébral | Par La rédaction avec AFP | 30/01/2024 |
| 38 | Un bébé de moins d'une heure retrouvé vivant dans un sac de courses à Londres | Par rédaction avec AFP | 19/01/2024 |
| 39 | Argentine : un incendie détruit 600 hectares d'un site Unesco | Par Belga | 27/01/2024 |
| 40 | EquiGala : Wilm Vermeir élu cavalier de l'année | Par Louis Lamote | 16/01/2024 |
| 41 | Philippe Henry (Ecolo) : un nouveau contrat de gestion pour les transports en commun, en plein déploiement en Wallonie | Par Par Olivier Arendt, d'après une interview de Thomas Gadisseux via La Première | 18/01/2024 |
| 42 | John Heymans pulvérise le record de Belgique du 5000m indoor et se qualifie pour les Jeux | Par Belga (édité par Alice Devilez) | 27/01/2024 |
| 43 | Les régimes militaires du Burkina, Mali et Niger décident de se retirer de la Cedeao | Par La rédaction Info avec AFP | 28/01/2024 |
| 44 | Le module lunaire japonais a repris vie, les analyses scientifiques vont pouvoir commencer | Par RTBF avec AFP | 28/01/2024 |
| 45 | Une personne décédée lors d'une attaque contre une église catholique italienne à Istanbul | Par La rédaction Info avec AFP | 28/01/2024 |
| 46 | Tensions entre le Pakistan et l'Iran : un problème local aiguisé par le climat régional | Par Pascal Bustamante | 18/01/2024 |
| 47 | Brexit : fin du blocage politique en vue en Irlande du Nord, après deux ans de paralysie | Par la rédaction avec Belga | 30/01/2024 |
| 48 | David Guetta et Swedish House Mafia enflammeront Tomorrowland 2024 | Par Belga avec RTBF Culture | 26/01/2024 |
| 49 | Atteinte aux droits d'auteur : le New York Times attaque en justice OpenAI, l'entreprise créatrice de Chat GPT | Par AFP | 28/12/2023 |
| 50 | Plusieurs actions menées par la police à Yser pour limiter le trafic de stupéfiants | Par Belga | 30/01/2024 |

Table 6: RTBF articles

## 7.2. Other IAA confusion matrices

Tables 7 and 8 represent the inter-annotator agreement matrices between annotators A and C, and B and C respectively. Agreement between A and B was already shown in Table 3. Annotator B was the main annotator.

| annotator A → vs C ↓ | English word/phrase | grammar mistake | longer piece of English text | non-linguistic remark | other linguistic remark | spelling mistake | strange/wrong construction | strangely/wrongly used word/phrase | word/phrase does not exist | #NA | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|
| English word/phrase | 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 14 |
| grammar mistake | 0 | 21 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 21 |
| longer piece of English text | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| non-linguistic remark | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 2 | 4 |
| other linguistic remark | 0 | 2 | 0 | 2 | 5 | 0 | 0 | 0 | 0 | 1 | 10 |
| spelling mistake | 0 | 0 | 0 | 0 | 0 | 11 | 0 | 0 | 0 | 5 | 16 |
| strange/wrong construction | 0 | 1 | 0 | 0 | 0 | 0 | 25 | 0 | 0 | 5 | 31 |
| strangely/wrongly used word/phrase | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 29 | 0 | 4 | 34 |
| word/phrase does not exist | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 16 | 0 | 16 |
| #NA | 1 | 2 | 0 | 1 | 0 | 8 | 2 | 2 | 0 | 23 | 39 |
| Total | 15 | 26 | 2 | 5 | 5 | 19 | 28 | 31 | 16 | 40 | 187 |

Table 7: Confusion matrix between annotators A and C

| annotator B → vs C ↓ | English word/phrase | grammar mistake | longer piece of English text | non-linguistic remark | other linguistic remark | spelling mistake | strange/wrong construction | strangely/wrongly used word/phrase | word/phrase does not exist | #NA | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|
| English word/phrase | 13 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 14 |
| grammar mistake | 0 | 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 21 |
| longer piece of English text | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| non-linguistic remark | 0 | 0 | 0 | 1 | 3 | 0 | 0 | 0 | 0 | 0 | 4 |
| other linguistic remark | 0 | 2 | 0 | 0 | 8 | 0 | 0 | 0 | 0 | 0 | 10 |
| spelling mistake | 0 | 1 | 0 | 0 | 0 | 12 | 0 | 0 | 1 | 2 | 16 |
| strange/wrong construction | 0 | 3 | 0 | 0 | 0 | 0 | 24 | 2 | 0 | 2 | 31 |
| strangely/wrongly used word/phrase | 0 | 2 | 0 | 0 | 0 | 0 | 1 | 27 | 1 | 3 | 34 |
| word/phrase does not exist | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 14 | 0 | 16 |
| #NA | | 10 | 0 | 0 | 2 | 5 | 9 | 4 | 1 | 8 | 39 |
| Total | 15 | 37 | 2 | 1 | 13 | 18 | 34 | 33 | 17 | 17 | 187 |

Table 8: Confusion matrix between annotators B and C

# Adding Argumentation into Human Evaluation of Long Document Abstractive Summarization: A Case Study on Legal Opinions

**Mohamed Elaraby**[†◇]**, Huihui Xu**[◇⋆]**, Morgan Gray**[◇⋆]**,**
**Kevin Ashley**[◇⋆]**, Diane Litman**[†◇⋆]

[†] Department of Computer Science, School of Computing and Information
[◇] Learning Research and Development Center
[⋆] Intelligent Systems Program, School of Computing and Information
University of Pittsburgh, Pittsburgh, PA USA
{mse30, hux16, mag454, ashley, dlitman}@pitt.edu

## Abstract

Human evaluation remains the gold standard for assessing abstractive summarization. However, current practices often prioritize constructing evaluation guidelines for fluency, coherence, and factual accuracy, overlooking other critical dimensions. In this paper, we investigate *argument coverage* in abstractive summarization by focusing on long legal opinions, where summaries must effectively encapsulate the document's argumentative nature. We introduce a set of human-evaluation guidelines to evaluate generated summaries based on argumentative coverage. These guidelines enable us to assess three distinct summarization models, studying the influence of including argument roles in summarization. Furthermore, we utilize these evaluation scores to benchmark automatic summarization metrics against argument coverage, providing insights into the effectiveness of automated evaluation methods.

**Keywords:** Summarization, Human Evaluation, Legal Summarization

## 1. Introduction

Human evaluation remains the best practice for evaluating generated summaries (Kryscinski et al., 2019; Fabbri et al., 2021), although conducting such evaluations can be laborious and costly, particularly when dealing with longform summaries exceeding 150 words (Krishna et al., 2023; Karpinska et al., 2021; Clark et al., 2021; Goyal et al., 2022b). Consequently, most longform summarization research shies away from conducting human evaluation (Krishna et al., 2023). While recent efforts have attempted to tackle this issue by standardizing the evaluation process with a focus on the factual accuracy dimension of the generated summaries (Krishna et al., 2023; Min et al., 2023) or coherence (Goyal et al., 2022b), none have adequately accounted for the unique requirements of the domain, which may entail additional dimensions.

In this paper, we propose the integration of a new dimension, **argument coverage**, into the human evaluation of abstractive summarization. We define *argument coverage* as the ability of the generated summary to adequately include argument components from the source document. Our focus lies on *long legal opinions*, a type of legal document mainly concerned with court decisions and characterized by intricate implicit argument structures dispersed throughout lengthy texts (greater than 4000 words on average) (Xu et al., 2021; Elaraby and Litman, 2022; Elaraby et al., 2023; Zhong and Litman, 2023). The summaries are mostly considered longform summaries (greater than 200 words

on average), Additionally, long legal opinions are composed of nuanced legal terminologies, necessitating legal experts for evaluation, which adds to the overall complexity of the task.

To address these research complexities, we make the following contributions: (1) We develop comprehensive human evaluation guidelines tailored for assessing argument coverage in generated abstractive summaries of long legal opinions. (2) We conduct a benchmarking study involving three existing systems, leveraging the introduced guidelines. This study aims to assess whether summarization models incorporating argument components achieve higher ratings of argument coverage compared to those that do not. (3) We assess the performance of automatic summarization metrics recently used in legal opinion summarization against human ratings, aiming to determine whether existing metrics adequately capture the variability in argument coverage within the generated summaries.

## 2. Related Work

Evaluating automatically generated summaries presents challenges such as scalability issues and low annotator agreement (Liu et al., 2023). These challenges are exacerbated when dealing with longform summaries, as assessing extended lengths inherently involves subjectivity (Karpinska et al., 2021). A comprehensive study by Krishna et al. (2023) revealed that 63% of research papers in longform summarization lack human evaluation. To

address this gap, they proposed guidelines for evaluating the factuality of longform summaries. Additionally, Min et al. (2023) introduced the FACTSCORE metric to assess the factuality of long-generated summaries (biographies), breaking down factuality into atomic facts for comparison against ground truth. Another framework by Chang et al. (2023) focuses on assessing coherence in book-length summaries by leveraging Large Language Model evaluation capabilities. *However, there is limited work addressing evaluation methods for legal documents, which often produce longform summaries.*

In the pursuit of evaluating generated legal summaries, Mullick et al. (2022) undertook a human assessment focusing on the relevance and readability of legal summaries. Similarly, Salaün et al. (2022) conducted a human evaluation to assess the fluency and adequacy of legal summaries. Xu and Ashley (2023) had a legal expert evaluator who indirectly assesses the information quality of legal summaries by evaluating the quality of generated question-answer pairs. *In this study, human evaluators directly evaluated the legal argument coverage in generated legal summaries.*

In efforts to benchmark automatic metrics against human evaluations, Fabbri et al. (2021) conducted a benchmarking study on automatic summaries generated from $23$ summarization models, sampled from the CNN-DailyMail dataset (Hermann et al., 2015). They evaluated these summaries using $14$ distinct automatic summarization metrics across dimensions such as factual consistency, coherence, fluency, and relevance. Building upon this work, Liu et al. (2023) expanded the evaluation framework to include Atomic Content Units (ACUs), which are fine-grained semantic units enabling high inter-annotator agreement. These new evaluation scores were used to augment benchmark summaries, including those from the news domain (CNN-DailyMail and Xsum (Narayan et al., 2018)) and the dialogue domain (SamSum (Gliwa et al., 2019)), against automatic metrics. *In our study, we focus on benchmarking automatic metrics used in legal opinion summarization against human evaluation scores for argument coverage.*

## 3. Dataset for Evaluation

In this analysis, we utilized a subset of the **CanLII** dataset [1], consisting of $1049$ cases annotated for argument roles types and summarization (Xu et al., 2021). The input legal opinions in this subset have mean and maximum lengths of $4375$ and $62786$ words, respectively, while the annotated summaries have mean and maximum lengths of $274$ and $2072$
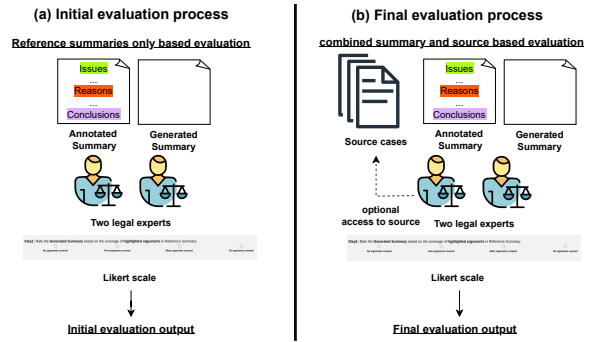


Figure 1: **Evaluation Process:** (a) *Initial evaluation* with human-annotated summaries and highlighted arguments. (b) *Final evaluation* with an option to cross-check the reference.

words, respectively. This subset has been extensively used in abstractive summarization research, particularly for constructing argument-aware abstractive summaries of legal opinions (Elaraby and Litman, 2022; Elaraby et al., 2023). The annotated argument roles follow the structure proposed in Xu et al. (2020, 2021), which breaks legal argument roles into three components: **Issue** (legal questions addressed by the court in the document), **Reason** (explanations for the court's decisions), and **Conclusion** (the court's rulings on the issues). Although these argument components constitute a small portion of the source cases, they typically account for $\approx 60\%$ of the summaries on average (Elaraby et al., 2023), highlighting the significance of considering argument roles in summary generation.

We considered the output of three different abstractive models in our evaluation process: $(1)$ **Finetuned LED-base:** This model serves as the baseline for legal opinion summarization, as described in Elaraby and Litman (2022). It finetunes the pretrained longformer-encoder-decoder (Beltagy et al., 2020) on the CanLII cases without additional information about the argument structure of the document. $(2)$ **arg-LED-base:** Utilizing the longformer encoder-decoder architecture, this model highlights argument units (Issues, Reasons, and Conclusions) with special tokens during both training and inference, as detailed in Elaraby and Litman (2022). $(3)$ **arg-aug-LED-base:** This model extends the arg-LED-base model, as discussed in Elaraby et al. (2023). It incorporates a mechanism for sampling summaries during inference and selecting the best model that exhibits the highest overlap with the input case's predicted argument roles.

---

[1]Data obtained through an agreement with CanLII (https://www.canlii.org/en/).

## 4. Argument Coverage Evaluation

We relied on two legal experts (two co-authors who are lawyers) to perform our human evaluation process, which was conducted in two phases. *Figure 1 shows an overview over the initial evaluation process (a) and the final evaluation process (b).*

### 4.1. Initial Evaluation Process

Initially, as shown in Figure 1 (a), we chose not to provide the full legal opinion due to its lengthy nature and the sparse distribution of argument roles across the case. Instead, experts were provided solely with human-written summaries, predominantly comprising argument roles. We highlighted the types of argument roles within the summaries to aid evaluators in distinguishing between argumentative and non-argumentative sections.

Our evaluation guidelines incorporate a 4-point Likert scale, facilitating a detailed assessment of argument coverage within the summaries. A rating of 4 indicates a perfect coverage of argument components, while a rating of 1 denotes a complete absence of coverage. To minimize misinterpretation of each score, we provided definitions for each rating category. During this phase, we utilized human-annotated summaries from 5 distinct legal opinions randomly selected from CanLII cases. For each case, we sampled summaries from the three distinct LED models, resulting in a total of 15 cases and summary pairs. Upon completion by both experts, the weighted quadratic kappa agreement, calculated using the sklearn implementation [2], between the two experts reached 0.466.

Discrepancies between the two experts were examined in a separate session, revealing that most disagreements stemmed from confusion regarding whether a certain argument within the generated summary was stated differently in the source document.

### 4.2. The Final Evaluation Process

To address evaluators' disagreements in the initial evaluation phase, we provided evaluators with human-written summaries, as outlined in the initial process. Additionally, evaluators were given the option to cross-check whether a specific argument was stated differently in the source document, as illustrated in Figure 1 (b).

Legal expert evaluators were provided with 15 additional summaries drawn from 5 new legal opinions. Our evaluation results suggest that by offering this option alongside the human-written summaries,

the overall weighted quadratic kappa agreement improved to 0.607. *The final evaluation guidelines are presented in Appendix A.*

### 4.3. Streamlining the Evaluation Process with Dedicated Software

To facilitate the experts' task, we developed a dedicated software for the longform evaluation of generated summaries. Our software builds upon the base code of the *Falte* tool (Goyal et al., 2022a), with several key enhancements: (1) **Keeping Expert State:** Recognizing the need for multiple sessions, we maintain the evaluation status for each expert, allowing them to complete the task across several sessions at their convenience. (2) **Inclusion of Likert Scale:** We include Likert scale definitions for each evaluation sample, aiming to reduce rating variability. (3) **Source Accessibility:** Acknowledging the positive impact of including source documents on the evaluation agreement, we added an option for experts to navigate to the source document. This allows them to cross-check confusing points against the source, improving accuracy. (4) **Highlighting Argument Roles:** To streamline the evaluation process, we highlight annotated argument roles in both the reference summaries and the source document. This facilitates cross-checking the generated summaries against them, reducing confusion. This approach is akin to solutions proposed by Krishna et al. (2023); Liu et al. (2023); Min et al. (2023), where evaluators are provided with atomic units of the summaries for evaluation. In our work, argument roles serve as the salient atomic units. *The tool is deployed and available online[3], enabling experts to complete tasks asynchronously. A screenshot is included in Appendix B [4].*

## 5. Results and Analysis

The final evaluation set consisted of 90 distinct generated summaries, that weren't included in the training phase, evenly selected from the three LED-based models, covering 30 unique legal opinion cases. Ratings were collected over two weeks using our dedicated software.

### 5.1. Experts' Agreement

The final quadratic kappa agreement was 0.483, which was lower than that obtained during the evaluation of the final evaluation process. *We hypothesize that this decline may be attributed to novel issues arising that were not addressed during the*

---

[2] https://scikit-learn.org/stable/modules/generated/sklearn.metrics.cohen_kappa_score.html

[3] https://summary-evaluation.herokuapp.com/

[4] https://github.com/EngSalem/legal-falte

| Metrics | $\tau$ correlation coeff. | | |
|---------|----------|----------|---------|
|         | **Expert 1** | **Expert 2** | **Average** |
| rouge-1 | 0.35 | 0.33 | 0.37 |
| rouge-2 | 0.33 | 0.30 | 0.33 |
| rouge-L | 0.28 | 0.34 | 0.34 |
| BERTscore | 0.31 | 0.29 | 0.33 |

Table 1: Automatic metrics correlations in *kendal tau $\tau$* with legal expert evaluations. All $\tau$ values are statistically significant with $p < 0.01$.

*training phase but required attention in the human guidelines.* We also evaluate the agreement between expert rankings of summaries by computing Kendall's tau ($\tau$) correlation coefficients. The final $\tau$ correlation coefficient is $0.429$ with $p < 0.001$, indicating a significant pairwise agreement between ratings of different systems.
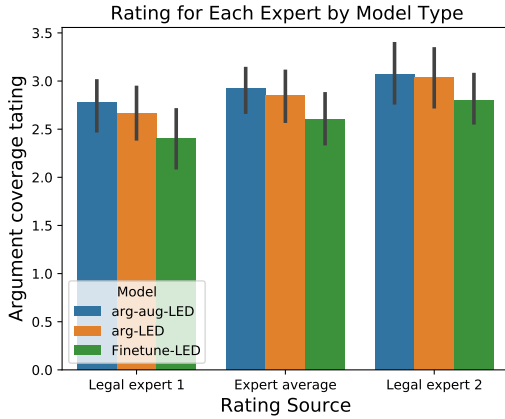


Figure 2: Average ratings. *Expert average:* average of Legal expert 1 and Legal expert 2.

## 5.2. Argument Aware Model Rankings

We analyzed the average rankings of summaries generated by different LED models. Figure 2 illustrates that the **Finetune-LED** model consistently received lower rankings from both legal experts compared to the **arg-LED** model (Elaraby and Litman, 2022), which highlights argument roles with special tokens, and the **arg-aug-LED** model (Elaraby et al., 2023), which leverages second-stage reranking to select the model with the highest argument similarity to the input. These findings are consistent with the significant correlation of rankings between both models discussed in 5.1, indicating that despite the drop in kappa agreement, experts agreed on the average rankings of summaries generated by different systems. *These results highlight that considering the argumentative components in the input document improves argument coverage in the generated summaries.*

## 5.3. Correlation with Automatic Metrics

We assess the effectiveness of automatic metrics previously employed in evaluating legal opinion summarization (Elaraby et al., 2023; Elaraby and Litman, 2022; Zhong and Litman, 2023) against human evaluation scores of argument coverage. These models primarily utilized ROUGE (Lin, 2004) and BERTScore (Zhang et al., 2019) to assess their proposed approaches. Table 1 shows that ROUGE demonstrated relatively higher correlation scores, ranging from $0.34$ to $0.37$, compared to BERTScore. Nevertheless, these findings suggest the potential for developing metrics specifically tailored to capture argument coverage. For instance, Fabbri et al. (2021) showed stronger correlations with aspects like fluency, consistency, coherency, and relevance, underscoring the need for more targeted metrics for assessing argument coverage.

## 5.4. Abstractiveness and Length of Summaries Effect on Ratings

Abstractiveness was quantified by computing the percentage of novel n-grams in each summary (See et al., 2017). Our findings, presented in Table 2, indicate that overall abstractiveness has limited influence on the ratings. However, as the number of novel n-grams increases (case of 4-gram), it can have a negative impact on argument coverage.

| Novel n-grams | Average | Expert 1 | Expert 2 |
|---------|---------|----------|----------|
| 1-gram | $-0.182^{\star}$ | $-0.151$ | $-0.180$ |
| 2-gram | $0.002$ | $0.001$ | $0.001$ |
| 3-gram | $-0.045$ | $-0.095$ | $0.002$ |
| 4-gram | $-0.200^{\star}$ | $-0.251^{\star}$ | $-0.129$ |

Table 2: $\tau$ values for novel n-grams vs ratings. $\star$ refers to $p < 0.05$.

Given the variability in our summary lengths, we aim to investigate its influence on argument coverage ratings. However, Table 3 indicates that the length of the summary has no significant effect on argument coverage.

| Expert Average | Expert 1 | Expert 2 |
|---------|----------|----------|
| 0.01 | 0.12 | $-0.08$ |

Table 3: $\tau$ values for summary length vs ratings. All values are with $p > 0.05$.

## 6. Conclusion

In this paper, we explored the concept of *argument coverage*, a new aspect in the evaluation of abstractive summarization. Our focus was primarily

on long legal opinions, where ensuring thorough argument coverage is essential for producing meaningful summaries. We introduced specific evaluation guidelines crafted for assessing argument coverage, allowing us to re-evaluate existing models for long legal opinion summarization. Our findings underscored the efficacy of integrating argument roles into the summarization process. Furthermore, we examined the automatic summarization metrics commonly used in legal opinion summarization research. Although ROUGE emerged as the most promising metric, our analysis suggests the potential for developing dedicated automatic metrics tailored to assess argument coverage more effectively. In future research, we aim to incorporate argument role types for a more nuanced evaluation and explore more efficient automatic metrics.

## Limitations

One limitation of this study is the absence of exploration into generated summaries from Large Language Models, which represents a promising avenue for future research in legal opinion summarization. Additionally, a larger dataset of legal opinions could have been incorporated into the evaluation training to refine the evaluation guidelines and potentially mitigate disagreements between experts more effectively. This would enhance the robustness of the evaluation process and bolster the reliability of the results. Moreover, while the focus was on legal opinions, extending the evaluation study to other domains where argument coverage is crucial, such as debates, would provide more comprehensive and inclusive guidelines for summarization.

## Acknowledgements

## 7. Bibliographical References

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

Yapei Chang, Kyle Lo, Tanya Goyal, and Mohit Iyyer. 2023. Booookscore: A systematic exploration of book-length summarization in the era of llms. *arXiv preprint arXiv:2310.00785*.

Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All that's 'human' is not gold: Evaluating human evaluation of generated text. 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 7282–7296, Online. Association for Computational Linguistics.

Mohamed Elaraby and Diane Litman. 2022. ArgLegalSumm: Improving abstractive summarization of legal documents with argument mining. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6187–6194, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Mohamed Elaraby, Yang Zhong, and Diane Litman. 2023. Towards argument-aware abstractive summarization of long legal opinions with summary reranking. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7601–7612, Toronto, Canada. Association for Computational Linguistics.

Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. SummEval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409.

Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*, pages 70–79, Hong Kong, China. Association for Computational Linguistics.

Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022a. FALTE: A toolkit for fine-grained annotation for long text evaluation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 351–358, Abu Dhabi, UAE. Association for Computational Linguistics.

Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022b. SNaC: Coherence error detection for narrative summarization. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 444–463, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. *Advances in neural information processing systems*, 28.

Marzena Karpinska, Nader Akoury, and Mohit Iyyer. 2021. The perils of using Mechanical Turk to evaluate open-ended text generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1265–1285, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Kalpesh Krishna, Erin Bransom, Bailey Kuehl, Mohit Iyyer, Pradeep Dasigi, Arman Cohan, and Kyle Lo. 2023. LongEval: Guidelines for human evaluation of faithfulness in long-form summarization. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1650–1669, Dubrovnik, Croatia. Association for Computational Linguistics.

Wojciech Kryscinski, Nitish Shirish Keskar, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Neural text summarization: A critical evaluation. 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 540–551, Hong Kong, China. Association for Computational Linguistics.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Yixin Liu, Alex Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. 2023. Revisiting the gold standard: Grounding summarization evaluation with robust human evaluation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4140–4170, Toronto, Canada. Association for Computational Linguistics.

Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore. Association for Computational Linguistics.

Ankan Mullick, Abhilash Nandy, Manav Kapadnis, Sohan Patnaik, Raghav R, and Roshni Kar. 2022. An evaluation framework for legal document summarization. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4747–4753, Marseille, France. European Language Resources Association.

Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.

Olivier Salaün, Aurore Troussel, Sylvain Longhais, Hannes Westermann, Philippe Langlais, and Karim Benyekhlef. 2022. Conditional abstractive summarization of court decisions for laymen and insights from human evaluation. In *Legal Knowledge and Information Systems*, pages 123–132. IOS Press.

Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.

Huihui Xu and Kevin Ashley. 2023. A question-answering approach to evaluating legal summaries. In *Legal Knowledge and Information Systems*, pages 293–298. IOS Press.

Huihui Xu, Jaromír Šavelka, and Kevin D Ashley. 2020. Using argument mining for legal text summarization. In *Legal Knowledge and Information Systems*, pages 184–193. IOS Press.

Huihui Xu, Jaromir Savelka, and Kevin D Ashley. 2021. Toward summarizing case decisions via extracting argument issues, reasons, and conclusions. In *Proceedings of the eighteenth international conference on artificial intelligence and law*, pages 250–254.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

Yang Zhong and Diane Litman. 2023. STRONG – structure controllable legal opinion summary generation. In *Findings of the Association for Computational Linguistics: IJCNLP-AACL 2023 (Findings)*, pages 431–448, Nusa Dua, Bali. Association for Computational Linguistics.

## A.   Final Evaluation Guidelines

Table 4 shows the final evaluation guidelines provided to legal experts to obtain argument coverage ratings.

## B.   Evaluation Tool

Figure 3 shows a snippet from the evaluation tool used for collecting argumentation coverage.

| Guide for Evaluation: Argument Coverage |
| :---: |
| **Description** |

*Argument Coverage: Do generated summaries cover the important points of the reference summary?*
You will be asked to rate the generated summary on a 4-point Likert scale to assess how well it covers the arguments in comparison to the highlighted arguments in the reference summary, which represent ground truth.

**Recommended Steps**

- Spend time to first read the reference summary until you understand the highlighted arguments.

- Read the generated summary until you understand its contents.

- Identify whether each argument highlighted in the reference summary is covered in the generated summary.

- If in doubt about a certain argument in the generated summary, click on the "go to source" button to double-check it against the source.

**Rating scale of the Generated Summary**

1. **No arguments covered:** The generated summary did not cover the highlighted arguments in the reference summary or covered them only inadequately.

2. **Few arguments covered:** The generated summary adequately covered only a limited number of the highlighted arguments in the reference summary.

3. **Most arguments covered:** The generated summary adequately covered most of the arguments highlighted in the reference summary.

4. **All arguments covered:** The generated summary adequately covered all the highlighted arguments in the reference summary.

Table 4: Final evaluation guidelines for argument coverage.



Figure 3: Screenshot from the tool used to collect argument coverage ratings from experts.

# A Gold Standard with Silver Linings: Scaling Up Annotation for Distinguishing Bosnian, Croatian, Montenegrin and Serbian

**Aleksandra Miletić[1]    Filip Miletić[2]**

[1]Department of Digital Humanities, University of Helsinki, Finland
[2]Institute for Natural Language Processing, University of Stuttgart, Germany
aleksandra.miletic@helsinki.fi    filip.miletic@ims.uni-stuttgart.de

## Abstract

Bosnian, Croatian, Montenegrin and Serbian are the official standard linguistic varieties in Bosnia and Herzegovina, Croatia, Montenegro, and Serbia, respectively. When these four countries were part of the former Yugoslavia, the varieties were considered to share a single linguistic standard. After the individual countries were established, the national standards emerged. Today, a central question about these varieties remains the following: How different are they from each other? How hard is it to distinguish them? While this has been addressed in NLP as part of the task on Distinguishing Between Similar Languages (DSL), little is known about human performance, making it difficult to contextualize system results. We tackle this question by reannotating the existing BCMS dataset for DSL with annotators from all target regions. We release a new gold standard, replacing the original single-annotator, single-label annotation by a multi-annotator, multi-label one, thus improving annotation reliability and explicitly coding the existence of ambiguous instances. We reassess a previously proposed DSL system on the new gold standard and establish the human upper bound on the task. Finally, we identify sources of annotation difficulties and provide linguistic insights into the BCMS dialect continuum, with multiple indicators highlighting an intermediate position of Bosnian and Montenegrin.

**Keywords:** BCMS, Distinguishing Between Similar Languages, human upper bound, gold standard, corpus annotation

## 1. Introduction

Bosnian, Croatian, Montenegrin and Serbian are the official standard linguistic varieties in their respective countries: Bosnia and Herzegovina (3.3M inhabitants), Croatia (3.9M), Montenegro (0.6M) and Serbia (6.7M) (Figure 1).[1] When the four countries were part of the former Yugoslavia, these varieties were considered to belong to the same language, which was commonly referred to as Serbo-Croatian or Croato-Serbian. After the civil wars of the 1990s and the establishment of individual countries, national linguistic standards also emerged. Thirty years later, one of the central questions about Bosnian, Croatian, Montenegrin and Serbian remains the following: How different are they from each other? In other words, how hard (or how easy) is it to distinguish between them?

One of the rare empirical studies that address this issue shows that Croatian and Serbian are situated at the opposing ends of the continuum, whereas Bosnian and Montenegrin tend to lean towards the one or the other depending on the considered linguistic feature (Ljubešić et al., 2018). Results from NLP, specifically on the task of Distinguishing Between Similar Languages (DSL) (Zampieri et al., 2014, 2017, 2015; Malmasi et al., 2016), seem to point in the same direction. In particular, Rupnik et al. (2023) introduce a four-class dataset for this task and evaluate two models. Model performance

varies widely per class: it is perfect on Serbian and solid on Croatian, but the results are weaker on Bosnian, and low on Montenegrin.

However, contextualizing model performance remains difficult since the human upper bound has not been determined. Furthermore, the four-class test set used in the system evaluation cited above allows only a single label per instance. Previous research has shown that this can be insufficient for DSL since some instances contain no variety-specific markers (Goutte et al., 2016; Bernier-Colborne et al., 2023; Zampieri et al., 2023). Finally, the dataset was annotated by a single human annotator. This may be suboptimal and potentially calls into question the reliability of the annotation, and thus of the evaluation.

This paper presents the first large-scale multi-annotator study on distinguishing Bosnian, Croatian, Montenegrin and Serbian (BCMS). Our goal is twofold. First, we seek to consolidate the existing four-class dataset by scaling up the number of annotators and introducing a multi-label annotation. Second, we systematically examine how human performance aligns with previous observations on the relationship between these varieties as well as system performance on the DSL task.

Our contributions are as follows. (1) We **release a new gold standard set with multiple labels per instance**[2] for the DSL task on BCMS, drawing on multiple annotations per instance and an annota-

---

[1]Note that the number of inhabitants is not directly equivalent to the number of speakers of each variety.

[2]https://doi.org/10.5281/zenodo.10998042

tor population originating from all target countries. (2) We use this dataset to **reassess a previously proposed computational system**, investigating performance differences with respect to the original single-annotator, single-label test set. (3) We **establish the human upper bound** on this task and identify sources of annotation difficulties. (4) We **provide linguistic insights into the BCMS dialect continuum**, with multiple indicators highlighting an intermediate position of the varieties spoken in Bosnia and Herzegovina and Montenegro. To the best of our knowledge, this is the first perception study on the BCMS language area. Moreover, our contributions underline the validity of our methodology for experiments based on human annotation, independently of the tasks and languages at hand.

This paper is organized as follows. We first summarize related work (§ 2), present our annotation procedure (§ 3), and introduce the resulting dataset (§ 4). We then examine it from three perspectives: reassessing an existing DSL system (§ 5), analyzing human accuracy (§ 6), and comparing human and system performance (§ 7). We conclude with a summary and outlook (§ 8).

## 2. Related Work

Empirical research into the relationship between Bosnian, Croatian, Serbian and Montenegrin remains scarce. To address this issue, Ljubešić et al. (2018) conduct a corpus-based dialectometric study. The authors look at the geographical distribution of 16 linguistic variables on phonological, morphosyntactic and lexical levels. The results situate Croatian and Serbian at the opposing ends of the continuum, whereas Bosnian and Montenegrin tend to align with the one or the other depending on the variable. Furthermore, the variables do not necessarily have an even spread over the continuum or the same frequency. For example, the opposition between ekavian and ijekavian forms (e.g. *dete* in ekavian vs. *dijete* in ijekavian, meaning 'child') is a distinguishing feature for Serbian (the only of the four national standards based on both the ekavian and the ijekavian pronunciation); it is also by far the most frequent feature identified in the corpus by Ljubešić et al. (2018). This asymmetry can be expected to make some varieties harder to identify.

This hypothesis is corroborated by current results in the DSL task on these varieties. A DSL shared task has been organized regularly by the VarDial Workshops since 2014, and Bosnian, Croatian, Montenegrin and Serbian have been part of it from the very first iteration, albeit as a three-class problem focusing on Bosnian, Croatian and Serbian (Zampieri et al., 2014, 2015; Malmasi et al., 2016; Zampieri et al., 2017). In more recent work, Rupnik et al. (2023) introduce a novel benchmark,



Figure 1: Bosnia and Herzegovina, Croatia, Montenegro, Serbia and neighbouring countries. Map data ©2023 GeoBasis-DE/BKG (©2009), Google.

containing two datasets: SETimes, based on newspaper texts in Bosnian, Croatian and Serbian; and a Twitter dataset containing instances in Bosnian, Croatian, Montenegrin and Serbian. Their evaluation of two DSL systems on the four-class Twitter dataset yields a global micro F1 score of 0.87 for both models, but the results vary widely per class: both models achieve 1.0 micro F1 on Serbian and 0.83 on Croatian, but the scores are somewhat lower on Bosnian (0.75–0.80) and drop significantly on Montenegrin (0.10–0.36).

Goutte et al. (2016) identify similar areas of difficulty. In their comprehensive overview of available DSL methods, the authors report that classifiers show a high degree of confusion when discriminating between Bosnian and Croatian texts. They call on six native speakers from the three countries to manually annotate the 12 most problematic instances, obtaining mean annotator accuracy of 16.6%. Some of the instances receive no correct annotations; in the gold standard, these are systematically labelled either as Croatian or as Bosnian.

These observations have contributed to a drive to redefine the DSL task. To this end, Zampieri et al. (2023) introduce a True Labels dataset for English, Spanish and Portuguese, which introduces the `both/neither` label for instances without any variety-specific markers. The dataset is annotated manually. In a similar vein, Bernier-Colborne et al. (2023) argue for framing DSL as a multi-label classification task and introduce such a dataset for four varieties of French. A model trained and evaluated on their dataset achieves an absolute gain of 0.225 on the macro F1 score on ambiguous texts.

The findings discussed above coalesce around two main points. First, the current four-class BCMS dataset would benefit from redefining the anno-

tation as multi-label. Also, to ensure annotation reliability and determine the human upper bound, the number of human annotators should be scaled up (the current version is annotated by a single annotator). Second, a systematic examination of human performance is required in order to better understand both the relationship between these four varieties and the issues faced by NLP systems.

We address these challenges as follows. We recruit 33 annotators from the four target countries and reannotate the test set from the dataset proposed by Rupnik et al. (2023). The collected annotation is used to derive a new, multi-label gold standard test set, against which we reevaluate an existing system. We measure inter-annotator agreement and determine the human upper bound on the task, thereby enabling a better contextualization of NLP system performance. Finally, we contrast human and system performance and draw conclusions about sources of difficulty and the underlying properties of the dialect continuum.

## 3. Annotation Process

This section describes our data collection. First, we present the original dataset on which this work is based (§ 3.1). Next, we provide details on the task definition (§ 3.2) and data preprocessing (§ 3.3). Finally, we describe the demographic structure of our annotator pool (§ 3.4).

### 3.1. Original Gold Standard

The original dataset was collected from the social media platform Twitter (rebranded as X in 2023) using the TweetGeo (Ljubešić et al., 2016) and Tweet-CaT (Ljubešić et al., 2014) tools. It has been part of the VarDial shared task on DSL since 2016 (Malmasi et al., 2016) as an out-of-domain test set for systems trained on the newspaper-based SETimes dataset. The current version of the dataset was published by Rupnik et al. (2023).

In the dataset, a single instance corresponds to the concatenation of all tweets produced by a given user. The dataset contains 614 instances (4,456,087 tokens) in total, with a strong skew towards Serbian (Table 1). Results obtained on this dataset may therefore be less reliable for the other three varieties. The instances contain 7,257 tokens on average. Occasional tweets in languages other than BCMS were not filtered out. The instances were labeled manually by a single annotator.

The dataset is split into train, dev and test in a 3:1:1 ratio. We conduct our annotation on the test set, allowing us to reevaluate a previously proposed DSL system, establish the human upper bound, and more generally assess the relevance of multi-label annotation for this task.

| | Label | | | | |
|---|---|---|---|---|---|
| Split | bs | hr | me | sr | Total |
| train | 45 | 53 | 34 | 236 | 368 |
| dev | 15 | 18 | 11 | 79 | 123 |
| test | 15 | 18 | 11 | 79 | 123 |
| Total | 75 | 89 | 56 | 394 | 614 |

Table 1: Label distribution in the original gold standard across data splits. **sr** = Serbia, **hr** = Croatia, **bs** = Bosnia and Herzegovina, **me** = Montenegro.

### 3.2. Task Definition

The basic task in our annotation process is defined as follows: for a given instance, determine the country you think the author is from. We explicitly avoid asking the participants to identify the language of the author, since the interplay between national, ethnic and linguistic identity in this language area is complex (see e.g. Ljubešić et al., 2018). A speaker living in country A may exhibit linguistic features consistent with variety A, but self-identify as speaking variety B, C or D based on their ethnic identity. Since we are interested in the geographic spread of linguistic features independently of perceived ethnic identity, we ask for the country of origin to limit this type of bias. This is also reflective of the model we reevaluate: it was trained on top-level web domains of each country, which correspond more closely to geographic origin than to language.

Participants can provide a two-level annotation. In case of ambiguity, they are instructed to choose the country they find the most appropriate as the first-level choice, and can add multiple optional annotations as their second-level choice. This is in line with the previously discussed recent developments of VarDial DSL-TL (discriminating between similar languages – true labels) datasets for English, Spanish and Portuguese, which introduce the `neither/both` label for instances without variety-specific linguistic markers (Zampieri et al., 2023). However, the instances in these datasets are much shorter, spanning several sentences instead of hundreds of tweets per instance in our case. It is therefore much less probable to find a fully ambiguous instance in the BCMS dataset. We still include the multiple choice option, both for its linguistic relevance and to estimate annotator uncertainty.

Participants are also asked to highlight text segments on which they based their decision. They can choose between two types of segments: linguistic indicators and world knowledge. Annotation guidelines illustrate linguistic indicators with phonetic, morphological, lexical and syntactic phenomena; world knowledge pertains to country-specific named entities (TV channels, political parties, cities etc.). The guidelines explain the difference between the two types of indicators and ask for deci-

sions not to be based solely on world knowledge.

Finally, annotators are asked to mark the spot in the instance where they reached their decision. They may also report offensive content through the interface. Annotation is run using `potato` (Pei et al., 2022); a screenshot is provided in Figure 2.

### 3.3. Data Preparation

Unlike the original manual annotation, which was based on unaltered tweet content, we preprocess the data. We remove retweets (reposts of another user's tweet) since they are not produced by the users themselves. We also anonymize URLs and mentions in tweets by respectively replacing them with `[link]` and `@ime` (meaning 'name' in BCMS). This is done for two main reasons: to avoid priming the participants based on the content of these elements, and to improve readability for participants not familiar with Twitter. Hashtags are left unaltered, since they are often part of sentence structure. A brief description of these elements and their processing was provided in the annotation manual.

### 3.4. Participants

Participants were recruited through the authors' personal and professional contacts. Participants needed to be legally of age, to be native speakers of one of the four varieties, and to have spent most of their lives in one of the four countries.

Participation was not remunerated. This fact, as well as the expected duration of the task, was clearly stated both in the call for participation and the informed consent form. Prospective annotators were required to email the authors, read the task instructions and informed consent form, sign it and return it by email. Their willingness to complete this process was taken as an indicator of their motivation to participate despite the lack of remuneration. Further details are provided in the Ethics Statement (Section 10).

A total of 33 participants were recruited. A pre-annotation survey asked for participants' gender, year of birth, place of birth, current country of residence, the country in which they spent most of their lives until now, and until the age of 18.

A total of 25 participants identified as female, and 8 as male. Mean annotator age was 44.6 (SD = 12.1). In the analyses presented here, we consider the participants to come from the country in which they spent most of their lives according to the pre-annotation survey. The distribution of participants per country is given in Table 2.

Note that not all of the participants annotated the full dataset. Because participation was not remunerated, we aimed to limit the expected task duration to 1h. To this end, we split the dataset into four subsamples. The first subsample had the

highest number of participants (17) and they were the most diverse. For the remaining three, most if not all participants were from Serbia.

| Country | Total | S1 | S2 | S3 | S4 |
|---|---|---|---|---|---|
| Bosnia and Herz. | 4 | 4 | — | — | — |
| Croatia | 7 | 6 | 1 | — | — |
| Kosovo | 1 | — | 1 | — | — |
| Montenegro | 1 | 1 | — | — | — |
| Serbia | 20 | 6 | 4 | 5 | 5 |
| Total | 33 | 17 | 6 | 5 | 5 |

Table 2: Distribution of participants by self-reported country. S1-S4: subsamples 1-4.

## 4. Establishing the New Gold Standard

The collected annotations were used to establish a new, multi-label gold standard. We describe how the new gold labels were determined (§ 4.1) and analyze the resulting label distribution (§ 4.2).

### 4.1. Resolving Annotations

Data collection ran from June to September 2023. After excluding participants who annotated less than 5 instances, the collected data contains a total of 1,098 annotations, out of which 988 were first-level annotations, and 110 were optional second-level choices. The median number of annotators per instance was 5 (min = 3, max = 17).

Inter-annotator agreement is evaluated using Krippendorff's $\alpha$ (Krippendorff, 1970), computed via the `Fast Krippendorff` implementation (Castro, 2017). As shown in Table 3, there are notable differences between the four subsamples, with $\alpha$ ranging from 0.668 on Subsample 3 to 0.893 on Subsample 1. This may be an indicator of sample difficulty, but further investigation is required to confirm this. All scores correspond to acceptable levels of agreement (Krippendorff, 2004).

We establish the new gold standard using a weighted voting strategy. The label selected as the first-level choice receives the weight of 1, and all

| Subsample | $\alpha$ |
|---|---|
| S1 | 0.893 |
| S2 | 0.734 |
| S3 | 0.668 |
| S4 | 0.768 |
| Average | 0.765 |

Table 3: Inter-annotator agreement measured as Krippendorff's $\alpha$. $-1 \leq \alpha < 0$: inverse agreement; $\alpha = 0$: no agreement beyond chance; $0 < \alpha \leq 1$: agreement beyond chance.

Figure 2: Annotation interface.

second-level choices receive the weight of 0.5. The votes are summed instance-level for each country and normalized by number of participants. Each country receives a final score between 0 and 1.

The gold first-level label is the one with the highest score. At this level, we do not accept multiple labels. One instance in the dataset did not receive a first-level annotation due to a tie in label scores and was excluded from the subsequent analyses. For the second-level annotation, we set a threshold at 0.2 in order to filter out labels which received low scores. In case of a tie on the second level, all labels with the second-best score are retained.

## 4.2. New Gold

The final label distribution in the new gold standard is given in Table 4. In the resulting annotation, 25 instances (20.3% of the dataset) have more than one label. For the instances that carry two labels, all combinations of countries are instantiated, except for the one combining Croatia and Montenegro. Note, however, that one instance in the dataset carries all four labels. This is also the only instance that has more than two labels.

| Label combo | Count | Labels | 1st | 2nd |
|---|---|---|---|---|
| sr | 70 | sr | 81 | 5 |
| hr | 16 | hr | 18 | 5 |
| bs | 7 | bs | 13 | 8 |
| hr, sr | 6 | me | 10 | 9 |
| bs, me | 5 | Total | 122 | 27 |
| me, sr | 5 | | | |
| bs, hr | 4 | | | |
| me | 4 | | | |
| bs, sr | 4 | | | |
| bs, hr, me, sr | 1 | | | |
| Total | 122 | | | |

Table 4: Distribution of labels in the new gold standard. Left panel: counts for all label combinations found in the new gold. Within a combination, labels are ordered alphabetically. Right panel: counts for each label as the first- and second-level choice.

Whereas Montenegro is the least frequent first-level annotation, it is the most frequent second-level choice (on 9 instances), followed by Bosnia and Herzegovina (on 8 instances). With Serbia and Croatia receiving only 5 second-level annotations each, this may point towards an uncertainty when it comes to identifying varieties from Bosnia and Herzegovina and Montenegro. This trend is explored in more detail in Section 6.

When compared to the original gold standard annotation, first-level labels differ on three instances. Two instances originally labelled as Montenegrin were relabelled as Serbian, and one instance initially annotated as Bosnian was recoded as Montenegrin. Such a low number of differences may be perceived as wasted annotation effort. However, the value of reliable annotations should not be underestimated. Moreover, the reannotation process had another goal: establishing a multi-label gold standard. This goal was achieved and its impact is evaluated in the following section. Finally, this process also allowed us to collect rich information on how humans perform on this task, which provide valuable observations laid out in Sections 6 and 7. We consider these as the silver linings of our work on the gold standard.[3]

## 5. System Evaluation

We examine the effect of changes to the gold standard on evaluations of DSL models. Specifically, we reevaluate the *NB Web* model introduced by Rupnik et al. (2023), which was the most robust in their evaluation. This is a Naive Bayes classifier trained on a web-based corpus using around 800 regionally distinctive words as features.

We compute the accuracy, macro-averaged and micro-averaged F1 scores using (i) the initial gold standard test set published by model authors; (ii) our reannotated test set in the single-label version; and (iii) a permissive evaluation, where a prediction is deemed correct if it corresponds to any

---

[3]This is also indicative of the reliability of the original annotator.

one label included in the multi-label version of our test set. The results are presented in Table 5.

| Gold standard | Acc. | F1 macro | F1 micro |
|---|---|---|---|
| initial | 86.9 | 67.7 | 86.9 |
| ours (one label) | 88.5 | 69.0 | 88.5 |
| ours (all labels) | 91.0 | — | — |

Table 5: Reevaluation of the DSL system by Rupnik et al. (2023). For comparability, initial test set results are recalculated to account for one instance excluded after reannotation.

The reannotated test set leads to a higher assessment of performance in the single-label setup (+1.6 accuracy points). Considering any label from the multi-label set as correct yields a further improvement (+4.1 accuracy points over the initial test set). These differences are overall limited – unsurprisingly, given the previously noted similarity between the initial and reannotated test sets – but they still confirm the relevance of multi-annotator and multi-label judgments on this task.

## 6. Human Performance

This section presents an analysis aiming to establish sources of difficulty for human annotators. We accomplish this by looking at two main indicators: annotators' accuracy as measured against the new gold annotation (§ 6.1), and their uncertainty (§ 6.2). For the latter, we rely on two indirect indicators: the presence of secondary labels and the duration of reading before the annotation decision is reached.

### 6.1. Accuracy

Compared against our single-label gold standard, mean participant accuracy on this task stands at 94.3 (SD = 6.2), or 5.8 points above model performance. It ranges from 76.7 to 100.0, indicating a considerable degree of variability across speakers.

To better understand the potential sources of this variability, we consider the available demographic information. We first check the effect of age under the assumption that older speakers may be better at distinguishing the varieties due to a higher degree of exposure prior to the breakup of Yugoslavia, but we find no correlation with annotator-level accuracy ($rho = -0.01$, $p = 0.97$).

We further look into the effect of the annotators' country of origin in relation to their accuracy on individual classes (Table 6). The analysis points to some intuitive patterns: for instance, speakers from Bosnia and Herzegovina obtain higher accuracy on instances labeled as coming from their country (+6.0) or from Croatia (+2.6) compared to speakers from Serbia, who are likely more susceptible to confusing those two varieties due to their shared

| Country | Accuracy on gold labels | | | |
|---|---|---|---|---|
| | bs | hr | me | sr |
| Bosnia & Herz. | 91.7 | 94.7 | 50.0 | 98.6 |
| Croatia | 96.4 | 100.0 | 93.3 | 100.0 |
| Kosovo | 75.0 | 80.0 | 100.0 | 100.0 |
| Montenegro | 75.0 | 100.0 | 100.0 | 100.0 |
| Serbia | 85.7 | 92.1 | 69.4 | 96.7 |
| Overall | 88.3 | 94.1 | 74.0 | 97.8 |

Table 6: Accuracy on individual gold labels cross-tabulated with annotators' self-reported countries of origin. Note that the number of annotators per country is highly variable.

frequent features (e.g. ijekavian forms). However, other patterns are less readily interpretable.

We further assess if class-level performance differs by country of origin using the Mann-Whitney–U test.[4] We compare the accuracy of annotators on a given class for one pair of countries at a time, and find no statistically significant differences.[5] The country-level trends may therefore be related to the uneven geographical distribution of annotators, but they should nevertheless be reexamined with a larger participant pool.

That said, Table 6 clearly shows that overall human performance varies strongly across the classes. Accuracy is highest on instances labeled as coming from Serbia and Croatia – the endpoints of the BCMS continuum – as opposed to those from Montenegro and Bosnia and Herzegovina. The Wilcoxon signed-rank test indicates that annotator accuracy is significantly different for all pairs of labels except Croatia and Serbia ($p = 0.129$). Variable degrees of difficulty in determining the correct label may also be reflected by other indicators of participants' uncertainty, to which we now turn.

### 6.2. Uncertainty

**Secondary labels.** Recall that participants annotated each instance using a primary country label and, optionally, one or more secondary labels. We now look into their tendency to use secondary labels as an indirect indicator of their uncertainty. Out of 988 individual annotations, 110 (11.1%) include a secondary country label. This tendency may seem overall limited; however, secondary labels were provided for 62 out of 123 annotated instances (50.4%). Furthermore, 28 out of 33 participants (84.8%) provided a secondary annotation at least once. This indicates that less-than-certain annotation decisions are in fact prominent.

---

[4] For all statistical significance tests, we set alpha to 0.05. Full results with individual test statistics and p-values are provided in Appendix A.

[5] We do not extend this analysis to the Kosovo and Montenegro groups as each only has one annotator.

| 1st choice annotations | | | 2nd choice labels | | | | Time to decision | | | Chars to decision | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Label | Total | w\ 2nd choice | bs | hr | me | sr | med. | min | max | med. | min | max |
| bs | 119 | 27 (22.7%) | — | 15 | 10 | 10 | 1' 32" | 0' 02" | 6' 04" | 1,200 | 0 | 5,239 |
| hr | 159 | 16 (10.1%) | 8 | — | 5 | 5 | 1' 12" | 0' 01" | 4' 52" | 1,028 | 105 | 5,173 |
| me | 68 | 19 (27.9%) | 10 | — | — | 12 | 1' 37" | 0' 01" | 7' 12" | 2,204 | 657 | 9,218 |
| sr | 642 | 48 (7.5%) | 28 | 10 | 23 | — | 1' 21" | 0' 01" | 5' 35" | 1,366 | 0 | 7,077 |
| Total | 988 | 110 (11.1%) | 46 | 25 | 38 | 27 | 1' 23" | 0' 01" | 7' 12" | 1,330 | 0 | 9,218 |
| (a) | | | (b) | | | | (c) | | | (d) | | |

Table 7: Distribution of individual annotations by choice of primary country labels. Panels from left: (a) number of annotations; (b) distribution of secondary country choices (may not sum to first country totals due to multiple choices being allowed); (c) time taken to annotate an instance; (d) character index where decision was made, indicated by highlighting tweet text. Outliers excluded in panels (c) and (d).

We further examine this trend with respect to different primary country choices under the assumption that different regional varieties are not equally easy to distinguish. The results in Table 7 show a clear distinction between annotations resulting in primary labels of Serbia or Croatia, with secondary choices present in up to 10% of cases; and those of Bosnia and Herzegovina or Montenegro, where secondary choices are two to three times more frequent. This is consistent with the intermediate position of these two countries in the regional dialect continuum (previously noted in Section 6.1).

The distribution of secondary labels varies depending on the primary country, but without clear tendencies: whatever the primary country choice, most (if not all) other countries may be considered as potential alternatives. These overlaps are striking as we would expect them to more clearly pattern with similarities between the varieties. We therefore conduct a qualitative analysis to better understand the motivations for secondary choices.

**Qualitative analysis.** Consider the following sample tweets (normalized to include diacritics), taken from a single instance where both primary and secondary choices hesitated between Montenegro and Bosnia and Herzegovina.

(1) Ako mi nestane interneta, umrijet ću.
*If I run out of internet, I will die.*

(2) Komšija pošalje poruku da mu lajkujem profilnu.
*A neighbor messaged me to like his profile pic.*

Example (1) includes the future tense form *umrijet ću* 'I will die', which is atypical for most of Serbia. It is the only dialect region where this construction would generally be realized with ekavian phonological features and fully synthetically (*umreću*). Example (2) contains the lexical item *komšija* 'neighbor'. Its use excludes Croatia, the one dialect region where the equivalent *susjed* is predominant. This would leave the annotator with the choice between the varieties of Bosnia and Herzegovina and Montenegro, which have many more shared linguistic features. In other words, the difficulty comes from insufficiently distinctive regional linguistic features.

A different pattern is illustrated by the following tweets, taken from an instance where annotators were hesitant between Montenegro and Serbia.

(3) Današnji dan – jedva čekam sjutra.
*Today – I can't wait for tomorrow.*

(4) i lep i jak
*both handsome and strong*

Example (3) contains the form *sjutra* 'tomorrow'. It distinguishes Montenegro from all other varieties, which have the equivalent *sutra*. But a minority of this user's tweets contain forms typical of varieties spoken in Serbia. Example (4) includes the ekavian variant *lep* 'pretty, handsome', whereas in Montenegro we would expect the ijekavian *lijep*. This can be seen as codeswitching. It is often spurious (e.g. quoted song lyrics), but codeswitched instances are not systematically flagged on Twitter. Annotation is therefore complicated by linguistic features which are sufficiently distinctive on their own, but which together point to multiple regional varieties.

**Duration of reading.** A final type of information on annotation difficulties comes from behavioral data: the automatically recorded amount of time spent to annotate an instance; and the character index at which the decision was made, indicated by highlighting tweet text. Distribution by primary country choice is shown in Table 7, panels (c) and (d). For each variable, we use the Mann-Whitney–U test to determine whether it differs significantly across individual pairs of labels.

Annotation duration varies depending on the chosen primary label. Annotators spend less time on instances they label as Croatian or Serbian, and more on those labeled as Bosnian or Montenegrin. The difference in median annotation duration is up to 25 seconds (Croatia vs Montenegro). These differences are statistically significant in all pairs of labels, except for those with a similar status in the dialect continuum: Croatia and Serbia, and Bosnia and Herzegovina and Montenegro.

Looking at the amount of read text, it is by far the highest when labelling an instance as coming from Montenegro – up to twice more compared to the other labels. The differences are statistically significant in all pairs of labels, except when comparing Bosnia and Herzegovina – which has the second lowest median – with Croatia and with Serbia. This is a slight reversal of the previous tendency; a potential explanation is that identifying features distinctive of Bosnia and Herzegovina requires somewhat less text, but more careful consideration, compared to those typical of Serbia.

Overall, behavioral information aligns with other indicators of annotation uncertainty: varieties at the extremes of the regional dialect continuum are easier to discriminate than those with an intermediate position. We now ask whether these trends also hold for system performance.

## 7. Human vs. System Performance

As previously noted (§ 6.1), mean human accuracy is noticeably higher than system performance on this task. We now compare human and system performance at a finer-grained level by contrasting their respective confusion matrices (Figure 3).



Figure 3: Confusion matrices for human and system performance. The matrix for annotators is computed on all individual annotations. The values are normalized per true label.

Both humans and the model obtain the highest results on instances labeled as coming from Serbia. The model in fact achieves perfect performance, potentially reflecting the skew in its training data.

The class with the second-highest human accuracy is Croatia, with the misclassified instances spread over all three remaining classes. The system obtains an accuracy that is over 10 points lower. Moreover, for misclassified items, it systematically falls back onto Bosnia and Herzegovina. We find a similar pattern for the class of Bosnia and Herzegovina: the system performs somewhat worse than our annotators and, unlike them, always misclassifies into the same class – in this case, Serbia.

Finally, both the annotators and the system struggle the most with instances labeled as coming from Montenegro, although to a very different extent. Our participants produce misclassifications in 26% of cases; in half of these annotations, they opt for Bosnia and Herzegovina, which again confirms the closeness of the two varieties. By contrast, the system misclassifies 90% of instances, splitting them between Bosnia and Herzegovina and Serbia.

## 8. Conclusions and Future Work

We have presented the first large-scale multi-annotator study on distinguishing Bosnian, Croatian, Montenegrin, and Serbian – four closely related but distinct national linguistic varieties. In order to consolidate an existing single-annotator, single-label test set for the task of Distinguishing Between Similar Languages, we scale up the number of annotators and recruit them from all target regions. This results in a multi-judgment, multi-label gold standard which allows us to analyze both system and human performance on this task.

Compared to the original test set, our reannotated version leads to a somewhat higher assessment of accuracy of an existing system (88.5, or +1.6 points, on single-label evaluation). More importantly, we establish mean human accuracy (94.3), showing that the system still lags behind it. We further identify sources of annotation difficulties using a broad range of indicators and observe consistent effects in line with the properties of the regional dialect continuum. These results may be partly due to an imbalanced geographic distribution of our annotators, but they point to important considerations which can be further validated on a larger participant sample. Specifically, instances coming from the endpoints of the dialect continuum – Croatia and Serbia – are the most accurately annotated and the easiest to judge; the reverse is true for Bosnia and Herzegovina and (especially) Montenegro, which occupy an intermediate position and have been shown to exhibit less distinctive features. Finally, a comparative error analysis shows that human misclassifications are spread across the false classes and likely explained by linguistic similarities. By contrast, the system generally falls back onto one dominant class, reflecting the label distribution in its training data.

Our results also raise questions to be explored in future work. The use of optional secondary labels in human annotation has shown that one-fifth of instances give rise to ambiguous interpretations. Formulating the DSL task as multilabel classification on these varieties would therefore more closely align model design with the perceptions of native speakers. More generally, the target varieties vary in terms of their relative annotation difficulty, with the one spoken in Montenegro proving particularly challenging. But this is also the most recently es-

tablished of the four national standards, suggesting an important role of diachronic developments. Additional annotators from as yet underrepresented countries would enable a further analysis of this and other empirically established patterns, providing novel insights into this linguistically rich region.

## 9. Limitations

A central aim of our study was to reannotate an existing dataset; we were therefore bound by its original class distribution. This however implies a strong skew towards data from Serbia, with Montenegro being the least frequent of the remaining three classes. This trend may have an impact on the analysis of human behavior, which could be verified through a replication study on a balanced subsample. A connected issue is the geographic skew in our annotator sample, as discussed throughout the paper. The reliability of the results is particularly affected for the Kosovo and Montenegro groups, with only one annotator each.

More generally, the annotated instances are long, with an average of over 7,000 tokens. Rather than request that annotators read all instances in their entirety – which does not seem reasonable in terms of cognitive effort – we asked them to take a decision as soon as they had seen sufficient linguistic indicators. We note that individual annotators differ with respect to the amount of text they deem necessary to read. In addition, this approach is not strictly comparable to computational models, which generally use all available text.

## 10. Ethics Statement

This study draws on data provided by 33 human annotators. All participants gave informed consent prior to accessing the annotation platform. The informed consent form described the task to be performed; the nature of the data to be annotated (tweets), including the risk of being exposed to potentially offensive content; the estimated duration of the task; the specific demographic information to be collected; the non-remunerated nature of participation; the right to withhold answers to any questions and to withdraw from the study at any moment; and the procedures used to anonymize and store the collected information. The participants could further freely opt into receiving the results of the study; being contacted for participation in extensions of the same study or in other similar studies; and being publicly acknowledged as participants in resulting scientific publications and dataset documentation.

We collected personal information on the participants: gender, year of birth, place of birth, and country-level residential history. We used this information to provide aggregate analyses of annotation performance and perception of different regional language varieties. Moreover, we aimed to fully respect the self-reported nature of this information. For example, in selecting their country of origin, one participant chose the option "other" and entered "Kosovo", while self-identifying as a speaker of Serbian in correspondence with the authors. We assigned this participant to the Kosovo group in line with their choice. Participant-level personal information is anonymized and securely stored. We disclose the names of a subset of participants in order to acknowledge their participation, but without linking the names to any other information. This was explicitly agreed through an opt-in procedure.

In terms of more general risks, we note that linguistic research in socially complex contexts – including areas with a history of conflict – may be instrumentalized with respect to broader societal or political issues. We stress that our research empirically examines regional patterns of language use as attested in the data we collected, without a predetermined view of the linguistic communities under study or suggestion that the observed patterns generalize to the population level.

## 11. Acknowledgements

## 12. Bibliographical References

Gabriel Bernier-Colborne, Cyril Goutte, and Serge Leger. 2023. Dialect and variant identification as a multi-label classification task: A proposal based on near-duplicate analysis. In *Tenth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2023)*, pages 142–151, Dubrovnik, Croatia. Association for Computational Linguistics.

Santiago Castro. 2017. Fast Krippendorff: Fast computation of Krippendorff's alpha agree-

ment measure. `https://github.com/pln-fing-udelar/fast-krippendorff`.

Cyril Goutte, Serge Léger, Shervin Malmasi, and Marcos Zampieri. 2016. Discriminating similar languages: Evaluations and explorations. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 1800–1807, Portorož, Slovenia. European Language Resources Association (ELRA).

Klaus Krippendorff. 1970. Estimating the reliability, systematic error and random error of interval data. *Educational and psychological measurement*, 30(1):61–70.

Klaus Krippendorff. 2004. *Content Analysis: An Introduction to Its Methodology*, second edition. SAGE publications.

Nikola Ljubešić, Darja Fišer, and Tomaž Erjavec. 2014. TweetCaT: a tool for building Twitter corpora of smaller languages. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 2279–2283, Reykjavik, Iceland. European Language Resources Association (ELRA).

Nikola Ljubešić, Maja Miličević Petrović, and Tanja Samardžić. 2018. Borders and boundaries in Bosnian, Croatian, Montenegrin and Serbian: Twitter data to the rescue. *Journal of Linguistic Geography*, 6(2):100–124.

Nikola Ljubešić, Tanja Samardžić, and Curdin Derungs. 2016. TweetGeo - a tool for collecting, processing and analysing geo-encoded linguistic data. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 3412–3421, Osaka, Japan. The COLING 2016 Organizing Committee.

Shervin Malmasi, Marcos Zampieri, Nikola Ljubešić, Preslav Nakov, Ahmed Ali, and Jörg Tiedemann. 2016. Discriminating between similar languages and Arabic dialect identification: A report on the third DSL shared task. In *Proceedings of the Third Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial3)*, pages 1–14, Osaka, Japan. The COLING 2016 Organizing Committee.

Jiaxin Pei, Aparna Ananthasubramaniam, Xingyao Wang, Naitian Zhou, Apostolos Dedeloudis, Jackson Sargent, and David Jurgens. 2022. POTATO: The portable text annotation tool. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 327–337, Abu Dhabi, UAE. Association for Computational Linguistics.

Peter Rupnik, Taja Kuzman, and Nikola Ljubešić. 2023. BENCHić-lang: A benchmark for discriminating between Bosnian, Croatian, Montenegrin and Serbian. In *Tenth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2023)*, pages 113–120, Dubrovnik, Croatia. Association for Computational Linguistics.

Marcos Zampieri, Shervin Malmasi, Nikola Ljubešić, Preslav Nakov, Ahmed Ali, Jörg Tiedemann, Yves Scherrer, and Noëmi Aepli. 2017. Findings of the VarDial evaluation campaign 2017. In *Proceedings of the Fourth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial)*, pages 1–15, Valencia, Spain. Association for Computational Linguistics.

Marcos Zampieri, Kai North, Tommi Jauhiainen, Mariano Felice, Neha Kumari, Nishant Nair, and Yash Bangera. 2023. Language variety identification with true labels. *arXiv preprint arXiv:2303.01490*.

Marcos Zampieri, Liling Tan, Nikola Ljubešić, and Jörg Tiedemann. 2014. A report on the DSL shared task 2014. In *Proceedings of the First Workshop on Applying NLP Tools to Similar Languages, Varieties and Dialects*, pages 58–67, Dublin, Ireland. Association for Computational Linguistics and Dublin City University.

Marcos Zampieri, Liling Tan, Nikola Ljubešić, Jörg Tiedemann, and Preslav Nakov. 2015. Overview of the DSL shared task 2015. In *Proceedings of the Joint Workshop on Language Technology for Closely Related Languages, Varieties and Dialects*, pages 1–9, Hissar, Bulgaria. Association for Computational Linguistics.

# A. Full statistical results for human behavior analysis

| Label | Countries of origin | | U | p |
|---|---|---|---|---|
| | Croatia | Bosnia & Herz. | 15.5 | 0.778 |
| bs | Serbia | Bosnia & Herz. | 32.0 | 0.502 |
| | Serbia | Croatia | 50.0 | 0.197 |
| | Croatia | Bosnia & Herz. | 17.5 | 0.257 |
| hr | Serbia | Bosnia & Herz. | 37.5 | 0.845 |
| | Serbia | Croatia | 49.0 | 0.117 |
| | Croatia | Bosnia & Herz. | 23.0 | 0.060 |
| me | Serbia | Bosnia & Herz. | 52.0 | 0.344 |
| | Serbia | Croatia | 44.0 | 0.111 |
| | Croatia | Bosnia & Herz. | 17.5 | 0.257 |
| sr | Serbia | Bosnia & Herz. | 35.5 | 0.700 |
| | Serbia | Croatia | 49.0 | 0.119 |

Table 8: Results of the Mann-Whitney–U test comparing annotator-level accuracy for a given gold label across pairs of annotators' countries of origin. The Kosovo and Montenegro groups are limited to one annotator each and are therefore not included in the analysis.

| Labels | | W | p |
|---|---|---|---|
| bs | hr | 26.0 | 0.053 |
| bs | me | 25.5 | 0.015 |
| bs | sr | 8.0 | 0.003 |
| hr | me | 10.0 | 0.001 |
| hr | sr | 16.0 | 0.129 |
| me | sr | 12.0 | 0.001 |

Table 9: Results of the paired Wilcoxon signed-rank test comparing annotator-level accuracy across pairs of gold labels.

| Labels | | time to annotate | | character offset | |
|---|---|---|---|---|---|
| | | U | p | U | p |
| bs | hr | 11001.0 | 0.007 | 3127.0 | 0.242 |
| bs | me | 3839.5 | 0.694 | 773.5 | 0.005 |
| bs | sr | 32842.5 | 0.046 | 13522.5 | 0.672 |
| hr | me | 6557.5 | 0.009 | 2115.0 | 0.000 |
| hr | sr | 53433.0 | 0.205 | 18759.5 | 0.041 |
| me | sr | 18419.5 | 0.046 | 4939.0 | 0.006 |

Table 10: Results of the Mann-Whitney–U test comparing instance-level behavioral information (time taken to annotate an instance; character offset at which the decision was taken) across pairs of gold labels.

# Insights of a Usability Study for KBQA Interactive Semantic Parsing: Generation Yields Benefits over Templates but External Validity Remains Challenging

**Ashley Lewis,[1] Lingbo Mo,[1] Marie-Catherine de Marneffe,[2]**
**Huan Sun,[1] Michael White[1]**

[1]The Ohio State University     [2]FNRS - UCLouvain

[1]{lewis.2799, mo.169, sun.397, white.1240}@osu.edu     [2]marie-catherine.demarneffe@uclouvain.be

## Abstract

We present our findings from a usability study of an interactive semantic parsing system for knowledge-based question answering (KBQA). The system is designed to help users access information within a knowledge base without having to know its query language. The system translates the user's question into the query language, retrieves an answer, then presents an English explanation of the process so that the user can make corrections if necessary. To our knowledge, our work is the most thorough usability study conducted for such a system and the only one that uses crowdworkers as participants to verify that the system is usable for average users. Our crowdworkers participate in KBQA dialogues using 4 versions of a system based on the framework by Mo et al. (2022) and answer surveys about their experiences. Some key takeaways from this work are: 1) we provide evidence for the benefits of interactivity in semantic parsing with human users and using generated questions in lieu of templated representations, 2) we identify limitations of simulations and provide contrasting evidence from actual system use, and 3) we provide an examination of crowdsourcing methodology, in particular the trade-offs of using crowdworkers vs. a specially trained group of evaluators.

**Keywords:** interactive semantic parsing, user study, usability study, crowdsourcing

## 1. Introduction

Semantic parsing endeavors to translate everyday language into formal representations of meaning, encompassing structures like API invocations, SQL, and SPARQL queries. Even as large language models (LLMs) make headway on the task (Xu et al., 2023), parsing still grapples with significant hurdles: (1) parsers struggle with accuracy, unable to reliably handle the inherent ambiguity and diverse paraphrases characteristic of natural language queries, and (2) the parsing process remains opaque to users, impeding their ability to comprehend and verify results (Chen et al., 2023).

Addressing these obstacles, recent work delves into *interactive semantic parsing*, which involves human users in providing feedback to enhance system accuracy. Mo et al. (2022) does this for knowledge-base question answering (KBQA), developing an interactive framework to correct parses that offers a step-by-step explanation of the predicted query and allows users to correct each step using natural language. In addition to this system, they release the INSPIRED (**IN**teractive **S**emantic **P**ars**I**ng for Cor**RE**ction with **D**ecomposition) dataset to facilitate further exploration of interactive semantic parsing for KBQA. This dataset uses target questions from the COMPLEXWEBQUESTIONS (CWQ) dataset (Talmor and Berant, 2018).

Using the INSPIRED dataset and the framework developed by Mo et al. (2022), we conduct a usabil-



Figure 1: Example of the user interface on Amazon Mechanical Turk, created using the ParlAI/Mephisto system (Miller et al., 2017; Urbanek and Ringshia, 2023).

ity study in which we have Amazon Mechanical Turk (henceforth MTurk) workers test four different versions of a system using an interface developed using the ParlAI/Mephisto platform (Miller et al., 2017; Urbanek and Ringshia, 2023), which allows dialogue tasks to be run on MTurk. An example of the
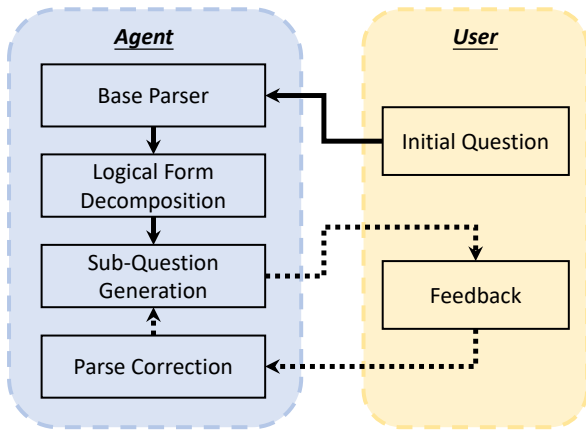
47

Figure 2: A high level overview of the framework described by Mo et al. (2022) for their interactive semantic parsing framework for KBQA.

interface can be seen in Figure 1, where the user is given the target question: *What is the official language of the country that contains [the] Al Sharqia Governorate?* For each dialogue, users are given a target question (originally from the CWQ dataset). These questions are multi-hop, meaning that they require more than one reasoning step to answer (hence the *complex* in COMPLEXWEBQUESTIONS). Mo et al. (2022) use a finetuned Transformer model (Vaswani et al., 2017) as a base automatic semantic parser to convert the question into a SPARQL query. They then decompose the query into pieces and translate those pieces to natural English sub-questions using a generation model.

We divert slightly from this method by first translating the query pieces into templated English using rule-based translation and then using a T5-based seq2seq model to translate the templates to natural language. This added step was implemented because query languages are unlikely to be well-represented in the pretraining data of the transformer, which makes the task more difficult (Kale and Rastogi, 2020). We then show these questions, along with intermediate answers retrieved from Freebase, to the user who can then use radio buttons to indicate whether or not the sub-questions and answers are correct. If they are not, the user can provide edits via further radio button selection (*insert, delete, replace,* or *edit*) and a text box for typing new questions (see Figure 1). The system then takes these edits and uses a parse correction model to try to generate a new parse, which again gets translated to natural-language sub-questions. This process, in theory, can repeat as many times as necessary. Figure 2 shows an overview of the system framework. To evaluate the viability of this system for real use, we recruited 48 crowdworkers to conduct dialogues in various conditions.

The contributions of this paper are as follows: 1) We provide a methodology for a thorough usability

study of the system using crowdworkers (Section 2). This methodology consists of a careful task design in which crowdworkers perform tasks in two systems and compare them. 2) We provide evidence of the benefits of the system designed by Mo et al. (2022) and identify areas for improvement (Section 3). In particular, we validate the choice to use interactivity, which had been shown to be helpful in previous work (Tian et al., 2023; Narechania et al., 2021; Yan et al., 2023; Elgohary et al., 2021), though mostly through user simulation. We provide external validity with human users. We also validate the use of a generation model instead of templated representations of questions. In terms of areas for improvement, we find that there are gains to be made in how dialogue context is best leveraged. 3) We examine the usability study with a critical eye and outline some lessons learned, which can help future research and usability studies (Section 4). We outline difficulties of Amazon Mechanical Turk and examine the trade-offs of using crowdworkers instead of a dedicated group of evaluators.

All data and models can be found on Github [1].

## 2. Usability Study Design

Our motivations for conducting a usability study were three-fold: 1) We wanted to verify that our system allows users of all types, even novices, to query knowledge bases with complex questions. The goal of our system is to elucidate the parsing process by decomposing the parse and translating it to understandable English sub-questions. It is designed to be intuitive and easy to use, so wanted to verify that this is indeed the case. 2) We wanted to identify problems and areas of improvement for the system. 3) We wanted to analyze how different aspects of the system influence user experience and success. We achieved this through ablations, testing a total of four systems.

### 2.1. System Designs

The first system is called the Full system, as it has no ablations. In this system, we use "fully contextualized" parse correction and generation models (see Figure 3 for further explanation). We use a model very similar to the best parse correction model from Mo et al. (2022), which is a finetuned T5 seq2seq model that takes the user's current correction, the original target question, and any previous sub-questions as input and outputs a corrected partial parse. The only difference in our model is that, instead of outputting a partial SPARQL query, it outputs a templated English version of the query. As mentioned, we do this because seq2seq models

---

Figure 3: A comparison of the Full and Decontextualized systems, which differ in the inputs on which the models are trained and tested. As you can see, the Full system is given more grounding information in the inputs of each model to use to produce an output. Example content can be found at the top of the diagram, which corresponds to the example in Figure 1. Not shown here is the Templated system which would follow the same method as the Full system but would terminate at step 3, as it does not use a generation model. The Non-Interactive system does not use either model.

tend to perform better on text-to-text tasks as opposed to text-to-query-language, as the query language is unlikely to be nearly as well represented in its training data (Kale and Rastogi, 2020).

Likewise, the generation model is also fully contextualized and translates templated English questions to natural language ones. Further, the first author conducted extensive experimentation to improve the generation model and thus we use the best version from those experiments. Full details can be found in Appendix A.

The remaining three systems each ablate an aspect of the Full system. The second system is a non-interactive baseline in which the user is simply shown the target question and its decomposition (natural language sub-questions and intermediate answers) and asked if it is correct. Users cannot make edits in this system. This allows us to gauge how informative the decomposition process is and establish a baseline by which to assess the value of interaction. While many papers do this via simulation (Mo et al., 2022; Yao et al., 2019; Yan et al., 2023; Elgohary et al., 2021), we wanted to know *in practice* how helpful it is.

The third system is decontextualized, meaning that neither the parse correction nor the generation models have access to the target question or previous sub-questions in their input. This means that they are mirrors of each other; the parse correction model takes a natural language sub-question (from the user) and produces a templated sub-question which can be translated to a partial SPARQL query using rule-based translation. The

generation model, meanwhile, takes that templated sub-question and outputs a natural language version. Mo et al. (2022) show that a contextualized model has better accuracy than a decontextualized one and thus in this study we seek to verify that this is the case. A comparison of the the Full and Decontextualized systems can be seen in Figure 3, which uses the same example as Figure 1.

The fourth system shows the templated sub-questions directly, not using the generation model at all. The other components remain the same as the Full system. We expected that the templates would be harder for users to understand in relation to the target question, based on the difference in BLEU scores between the natural language questions and templates in the INSPIRED dataset when referenced with the target question. The templates have an average of 2.13% for BLEU-1 scores, while the natural language questions score 17.59%. Together, the third and fourth systems allow us to assess the impact of the generation model and the third system allows us to investigate the role of context in both models.

To compare these systems, we set up three head-to-head comparisons of the Full system versus each of the ablated systems. A given crowdworker moves through a pipeline in which they complete 10 dialogues in one system (full or ablated), complete a 5-question survey about those dialogues, then complete 10 more dialogues in the other system and the survey about the second group of 10 dialogues, plus 3 questions that compare the two systems. At the end of each survey there is a com-

| Comparison | System | Officially Correct | Inferably Correct | User Deemed Correct |
|---|---|---|---|---|
| Full vs. Temp | Full | **57.75** | **61.75** | **74.25**[†] |
|  | Temp | 56.50 | 58.50 | 65.00 |
| Full vs. Decon | Full | **59.25** | **62.00** | 69.75 |
|  | Decon | 58.75 | 61.75 | **70.75** |
| Full vs. Non | Full | 59.00[†] | 62.00[†] | 70.75 |
|  | Non | 16.25 | 19.25 | - |

Table 1: Success rates of each system in each comparison (Full vs. Templated, Full vs. Decontextualized, and Full vs. Non-Interactive). † indicates a significant difference at a p-value of less than 0.05 using a paired-sample T-test.

ment box for any other feedback.

We estimated that workers would be able to complete one round of the comparison pipeline in 1 hour, estimating that each dialogue takes 2-3 minutes to complete and the surveys about 2 minutes apiece. We allow workers to do all three comparisons if their work is adequate, though workers never see a given target question more than once.

## 2.2. Item Selection

For this study, we wanted to balance selecting target questions that are representative of the INSPIRED test set, but also are exceptionally challenging, to thoroughly test and tease out the differences between systems. After conducting a power analysis with an estimated effect size of 0.2, we selected 400 questions that adequately represent all the categories and facets in the INSPIRED dataset (and thus the CWQ dataset on which it is based), but also over-represent characteristics that make queries more difficult for the system to predict. An overview of this selection can be found in Appendix B.

## 3. Results

The following sections outline the results of the usability study. In Section 3.1, we look at the success of the dialogues in each system and number of edits used. In Section 3.2, we compare the systems based on the surveys that workers completed.

## 3.1. Dialogue Results

To gauge the differences between the systems, we first look at the success rates of the 400 dialogues in each system. By *success*, we mean that the user is able to reach a correct final answer. This measure, however, is slightly more nuanced than it might first seem. In Table 1, we use three different measures of success. The first is *Officially Correct*, which means that the user confirmed that the sub-questions and answers were correct and the final answer exactly matched the gold answer for that query. The second is *Inferably Correct*, which

means that, when the user confirms that the decomposition is correct, the correct answer is easily inferable from the final answer.

(1) **Target Question:** When did the sports team that plays at Hilltop Park last win a championship?
**Gold Decomposition:**

1. What team plays at Hilltop Park?
   *Answer: San Francisco Giants*

2. What World Series has that team won?
   *Answer: 2010 World Series, 2012 World Series*

3. When did these events occur?
   *Answer: 2010 World Series: 2010, 2012 World Series: 2012*

4. Of these, which is the most recent?
   *Answer: 2012 World Series*

In (1), one can easily discern the final answer from the answer of sub-question 2. This is not always the case for superlative questions; often the final two questions are necessary. However, in cases where it is not, we count the dialogue as Inferably Correct.

The third category is *User Deemed Correct*, in which we determine that the user is satisfied with the answer they receive even if it is not correct. We include this category due to a curious trend we observe in which users seem to frequently misunderstand the target question. This is not a new phenomenon; Yao et al. (2019) observe the same issue in the human evaluation of their system. Our suspicions are strengthened by two factors: first, at the end of each dialogue and survey, users are given a text box to give feedback and flag any confusions or concerns, in which they frequently comment on the difficulty of understanding the target question. Second, we strongly emphasize in the tutorial and crowdworker communications that for each dialogue they can perform a maximum of 5 edits and that they must use all 5 if the decomposition and/or answer is incorrect. We find, however, that users frequently do *not* do this and in the majority of those cases, they have reached an answer to the dialogue (meaning that the final answer is not *None*). This seems to indicate that they believe the answer to be correct.

Thus, for the *User Deemed Correct* category, we include dialogues in which either 1) the user does not use all 5 of their available edits but reaches an answer even if it is incorrect, or 2) the user comments that they believe the decomposition to be correct but the final answer is *None*. We suspect that more questions could fall into this category, particularly because the user must answer affirmatively

the question *Are the above sub-questions accurate in relation to the complex question?* in order to end the dialogue. However, because we cannot verify the user's interpretations and it is possible that they answer *yes* to this confirmation question for reasons other than believing it to be accurate (such as wanting to move on to the next task), we conservatively define this category. We also omit a small number of cases in which there were generation model errors (18 in total), in case they misled users into thinking a decomposition was correct when it was not. Further explanation and examples of these omissions can be found in Appendix C. The issue of users misunderstanding the target question will be discussed further in Section 3.3.

Table 1 shows that the Non-Interactive system is significantly less successful on all measures than the Full system, demonstrating that interactivity greatly improves performance, which validates results seen in simulation by others (Elgohary et al., 2020, 2021; Yan et al., 2023; Yao et al., 2019). We also see that the Decontexualized system does not perform significantly differently than the Full system, perhaps indicating that context is not as important as we anticipated. Last, we find that the Templated system performs significantly worse than the Full system in the category of User Deemed Correct, which seems to indicate that users found it easier to extract meaning out of the generated sub-questions than the templated ones.

Figure 4 shows the breakdown of how many edits were used for each dialogue in the Templated and Full systems and the success rates (based on the *inferably correct* category). It shows that dialogues tend to be most successful after 1 or 2 edits, and success drops dramatically around 4 edits. The same graph for the Full and Decontextualized comparison, which had no significant differences, is in Appendix D. The comparison between the Full and Non-Interactive systems is not shown because the latter does not allow for edits.

Figure 4 also shows that there are many more 5-edit unsuccessful dialogues in the Templated system than the Full system. This seems to be due in large part to users stopping earlier with the Full system because the questions were easier to understand. Of the dialogues that have 4+ edits, we observe that 36.05% of them in the Templated system are cases where the user was able to get the correct final answer, but continued making edits, compared to 20.56% in the Full system. We assume this is because the templated questions were confusing enough to be misleading. It is noteworthy that our parse correction model was blocked from producing the same corrected question twice, except when the user used the exact same language in their feedback. This means that sometimes, particularly in the Templated system, the user did not

| Comparison | System | Effort Metric | Difference |
|---|---|---|---|
| Full vs Temp | Full | 0.2877 | 0.0519[†] |
| | Temp | 0.2358 | p = 0.0054 |
| Full vs Decon | Full | 0.2780 | 0.001498 |
| | Decon | 0.2795 | p = 0.9463 |

Table 2: Efficiency of dialogue systems. The table summarizes the effort metric comparisons between Full and Templated systems, and Full and Decontextualized systems. Statistical significance is indicated with †, which represents a p-value less than 0.05 using an approximate randomization test.

recognize that they had the correct answer and continued editing, unable to reproduce the correct answer again. See Appendix E for an example.

In evaluating the efficiency of the dialogue systems and assessing the relative worth of each edit made, we employed a metric of effort defined as the ratio of successful (inferably correct) dialogues to the total number of edits made, regardless of the dialogue's outcome. These numbers can be seen in Table 2. To rigorously test the significance of this observed difference, we employed approximate randomization testing. Details of this can be found in Appendix F.

The significant difference in effort between the Full and Templated systems strongly suggests that the type of system indeed impacts the efficiency of dialogues. The difference between the Full and Decontextualized system was non-significant, however, again suggesting that the systems did not perform substantially differently. This result seems to demonstrate the superior efficiency of the Full system over its Templated counterpart in facilitating successful dialogues with fewer edits.

This result, along with the *User Deemed Correct* difference seen in Table 1, provides evidence that the natural language questions are easier to understand than the templated questions. This tracks with our assumption that the natural language questions in INSPIRED more closely resemble the target question.

## 3.2. Survey Results

Table 3 shows the results of the surveys that each worker completed (2 per worker, per comparsison). The top section (questions 1-5) shows the average scores on the survey that workers completed after every batch of 10 dialogues in a given system. These results give us insights into user impressions of the different systems and we can see that they tend to cluster around the middle values. There are no significant differences between the Full and Decontextualized systems and only one significant result between the Full and Templated systems, for *I felt confident using the system*, in favor of the

Figure 4: A comparison of the success rates for the Full vs. Templated systems (using three different categories – officially correct, inferably correct, and user deemed correct) and the number of edits attempted by the user.

| | Question | Comparison 1 | | Comparison 2 | | Comparison 3 | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | Full | Temp | Full | Decon | Full | Non |
| 1 | I thought the system was easy to use. | 2.17 | 2.35 | 2.10 | 2.05 | 2.11 | **1.67**$^\dagger$ |
| 2 | I would imagine that most people would learn to use this system quickly. | 2.24 | 2.56 | 2.17 | 2.05 | 2.16 | **1.71**$^\dagger$ |
| 3 | I felt confident using the system. | **1.93**$^\dagger$ | 2.42 | 1.98 | 1.95 | 2.11 | 1.93 |
| 4 | The system was usually able to answer the questions correctly. | 2.43 | 2.49 | 2.24 | 2.19 | **2.32**$^\dagger$ | 3.00 |
| 5 | It took a reasonable amount of time to complete a dialogue. | 2.29 | 2.19 | 2.10 | 2.00 | 2.25 | **1.67**$^\dagger$ |
| 6 | Which of the two systems was easier to use? | 58.97 | 41.03 | 43.90 | 56.10 | 30.00 | **70.00**$^\dagger$ |
| 7 | Which system answered more questions correctly? | 56.41 | 43.59 | 53.66 | 46.34 | **67.50**$^\dagger$ | 32.50 |
| 8 | Which system gave you more confidence that the final answers were correct? | 58.97 | 41.03 | 43.90 | 56.10 | **70.00**$^\dagger$ | 30.00 |

Table 3: Survey results for each comparison. Users completed a survey of the first 5 questions after completing 10 dialogues in a given system. They answered using Likert scale responses: 1) Strongly Agree, 2) Agree, 3) Disagree, 4) Strongly Disagree (Lower scores are better). Significance was determined using a paired sample t-test. The last three questions were answered after completing a full comparison pipeline (10 dialogues in each system) and users selected only one of the two options for these questions. Significance was determined using a binomial test. $\dagger$ indicates a significant difference with a p-value of less than 0.05.

Full system. This indicates that the templated sub-questions increased users' confusion and eroded confidence. Unsurprisingly, there are significant results in the Full vs. Non-Interactive systems comparison. For the first two questions (*I though the sytem was easy to use* and *I would imagine that*

*most people would learn to use this system quickly*) and the final one (*It took a reasonable amount of time to complete a dialogue*), the difference is in favor of the Non-Interactive system, but this likely has a great deal to do with the lack of interaction. Users merely had to answer the question "Are these

52

sub-questions correct in relation to the Target Question?" and leave comments if necessary, making the task much easier and less time-consuming than the other systems. It is notable that question 4 had a significant result in favor of the Full system, again suggesting that interactivity greatly boosts system accuracy and user confidence in the final answer.

The bottom section of Table 3 (questions 6-8) shows the results of the comparative survey questions, which users answered after completing 10 dialogues in each system (Full vs. ablated system). Significant results only appear in the Full vs. Non-Interactive comparison for the reasons given above.

### 3.3. Discussion

This usability study revealed results that could be useful for semantic parsing research, in particular about the importance of interactivity, the benefit of generated questions over templated ones, and users' ability to successfully navigate such a system. However, there were fewer significant results than we anticipated, in particular between the Full and Decontextualized systems. This discrepancy, while somewhat disappointing, is elucidating in many respects. First, this demonstrates the difference between *simulation* and *real use*. Mo et al. (2022) use BART-large (Lewis et al., 2020) and QGG (Lan and Jiang, 2020) to simulate user feedback in their framework and find that they can get exact-match accuracy up to 73.5%, which is 14.25% higher than the best accuracy in our usability study (see Table 1). It is important to note, however, that their metric is calculated over the entire test set of INSPIRED, while ours is over a subset of 400 questions that were selected for their difficulty. They also find that using context in the model inputs leads to accuracy gains in simulation, a trend that did not bear out with human users in our study. Other works, such as (Elgohary et al., 2020, 2021; Yan et al., 2023; Yao et al., 2019), also show substantial gains using simulated users.

Our usability study shows how using a human user changes the overall success of the system. Our users, of course, are not *real* users, meaning that they were not using our system to answer their own questions, but rather questions given to them. There are two main reasons for this design choice; the first is that we wanted to directly test different systems on the same questions for the sake of a clear comparative analysis. Secondly, we wanted to lower the burden on the crowdworkers, who might have had difficulty coming up with their own complex questions that could be answered by the knowledge base.

One major factor to this difference in simulation and human users is the amount of noise that can be found in the dataset. While the INSPIRED dataset underwent a fair amount of cleaning to ensure that the

target questions matched their SPARQL queries, problems still remained. The questions, which came from the COMPLEXWEBQUESTIONS dataset (Talmor and Berant, 2018), are often very challenging to understand. They were created from the WEBQUESTIONSSP dataset (Yih et al., 2016) by combining two simple questions together with a rule-based method and then having crowdworkers rephrase them into a single, more natural-sounding question. This results in complex questions that are often quite unnatural and unlikely to be asked by a real human user, which in turn makes them more difficult to understand for our crowdworkers. We notice too a number of errors in the CWQ dataset – a mismatch between the target question and the meaning of the associated SPARQL query. This does not, however, seem to prevent gains in terms of parsing accuracy – systems such as Mo et al. (2022) and Niu et al. (2023) find that they are still able to make significant progress on correcting parses despite such errors.

In our study, however, users report frequently that the target questions were unnatural, incomprehensible, or unclear, demonstrating that such errors have a much bigger impact when using human users and in natural language generation settings. It also did not help that we chose the hardest questions available, which was intended to demonstrate that the system works well even in these cases. While this does appear to be the case the majority of the time, we underestimated users' fatigue and frustration with simply *understanding* the target question. Yao et al. (2019), as mentioned, notice similar problems in their user evaluation.

### 4. Lessons Learned: Crowdsourcing

The issue of the target questions appears to cause substantial problems in the overall accuracy of the systems, but also in terms of user investment in the task. We know based on worker feedback that these issues led to frustration and a slow-down of the tasks. We also hypothesized that workers' performance might improve over time, which did not occur. Looking at the first 10 and last 10 dialogues in the Full system of the 41 workers who did at least two comparisons, we find that the first 10 have an average success rate of 65.37% and the last ten have a rate of 62.44%.

Cultivating worker investment in a crowdsourcing task is a challenging but very important problem to consider. We carefully planned the crowdsourcing effort with several factors in mind, including ensuring that workers were compensated fairly and that the task took a reasonable amount of time to complete. We decided on a base rate for each dialogue according to the minimum wage of our state and the amount of time on average a dialogue might

take, estimated from internal testing. However, because users needed to read a tutorial and complete a qualification quiz prior to being granted access to our task, we also ensured that they were compensated for that time spent by granting them a bonus once they complete the whole pipeline, which we anticipated taking at most 1 hour. If a worker completed a full comparison they collected a total of $15. Because workers could complete all three comparison pipelines (assuming there were dialogues available), they could receive a maximum of $45. This payment rate is significantly higher than what is usually available on MTurk.

It is important to note, however, that using MTurk presents a great number of challenges. First, it is not really designed for dialogue tasks and external software is needed to do this. To this end, we utilized the ParlAI and Mephisto platforms (Miller et al., 2017; Urbanek and Ringshia, 2023), which are designed to run more complex tasks (particularly dialogue) than the regular MTurk platform. However, we find that the connection between the platforms can be very brittle and we encountered both a number of connectivity problems for our users.

Secondly, the issue of bots on MTurk poses a major challenge. Careful design went into preventing bots from accessing our task and eliminating ones that were able to get through, while not making the task too cumbersome for real crowdworkers. For example, we implemented a wait time for repeated attempts on the qualification quiz to prevent repeated random guessing on the questions until the correct answers were found. We also monitored how long it took a user to complete a task and manually reviewed tasks that were notably shorter than average. We also reviewed cases where there were no edits on several tasks in a row and where the same edit was submitted for every turn.

Further complicating this issue is that it is very important to *not* reject crowdworkers' completed tasks unless we are absolutely certain that it came from a bot. Due to the structure of the interaction between crowdworkers and requesters on MTurk, workers are unduly harmed by the rejection of work – it can irreparably damage their reputation on the platform and prevent them from being able to get future work. Thus, in cases of inadequate performance, if there is any doubt whether the worker is a bot or a human, we err on the side of caution and accept the work, but prevent that MTurk user from completing more tasks. We discarded the data from our analysis. In our experience, this problem appears to be worsening over time.

Thirdly, we suspect that because MTurk tasks are typically shorter and less involved, workers might have been less inclined to do our task in the first place and/or felt incentivized to go through the task

| Comparison | System | Officially Correct | Inferably Correct | User Deemed Correct |
|---|---|---|---|---|
| Full vs. Temp | Full | **72.2** | **77.8** | **81.5** |
| | Temp | 68.5 | 68.4 | 68.4 |
| Full vs Decon | Full | **63** | **65.2** | **67.4** |
| | Decon | 52.9 | 52.9 | 54.9 |
| Full vs. Non | Full | **63.0** | **68.5** | 76.0 |
| | Non | 18.0 | 24.5 | - |

Table 4: Percentage success rates of a smaller, possibly more dedicated, working group of 6 crowdworkers that participated in our task. These numbers represent success out of 60 items (10 per worker).

as quickly as possible, which might have resulted in a decrease in performance.

This method of human evaluation stands in contrast with a method that seems increasingly typical: hiring a dedicated, small team of evaluators and specifically training them to do the task. There are obvious benefits to this approach: it is much easier to ensure that your workers have understood the task, it is possible to conduct more complicated/involved tasks, and you can communicate directly and easily as problems arise. It is thus much easier to ensure quality results.

However, drawbacks of this include the difficulty of finding such workers, the added expense involved, and – importantly – the possible detraction from the generalizability of the results. We opted to use crowdworkers in large part because we wanted to ensure that the task was easily understandable and doable for an average person without a great deal of guidance. However, we find that this benefit seems to be getting washed out by the difficulties of conducting a crowdsourcing task on MTurk. This finding is supported by looking at the performance of a particular subset of our crowdworkers. On a previous crowdworking task, the first author established a relationship with a "working group" on MTurk. Such groups have organically formed on online spaces; workers collaborate to share information about tasks they can work on, requesters, etc. This particular group had produced excellent work in a previous task, so the first author reached out to them again and invited them to work on this one. Six of them agreed and did all three comparisons. By looking at their work in isolation, we see that the results look much more like what we expected. The success rates of the working group can be seen in Table 4. This suggests that a more dedicated group might have yielded clearer results than a general crowdsourcing task.

## 5.   Related Work

As mentioned, this usability study is, to our knowledge, the most thorough one conducted on an interactive semantic parsing system. There are a

number of systems developed similar to that of Mo et al. (2022) though mostly in the realm of text-to-SQL parsing instead of KBQA. This includes recent work such as Chen et al. (2023), who use CodeBERT as an error detector for SQL parsing tasks, which they envision being incorporated into an interactive system. Also included is Eyal et al. (2023) which uses a decomposition strategy for SQL queries, Stengel-Eskin and Van Durme (2023) who leverage confidence measures in an interactive SQL parsing system to help annotators make corrections, and Yan et al. (2023) who simulate natural language feedback.

In terms of KBQA, Xu et al. (2023) fine-tune LLaMA (Touvron et al., 2023) and pair it with GPT-3 to produce parses, while Niu et al. (2023) use a decomposition strategy to improve a parser's understanding of natural language. However, these systems stand in contrast to that of Mo et al. (2022) as neither of these systems are interactive.

Narechania et al. (2021), Ning et al. (2023), Yao et al. (2019), and Tian et al. (2023) all create interactive systems for text-to-SQL parsing and conduct user studies for their systems using in-house participants (students or employees from their institutions), many of which have experience with SQL. Additionally, Tian et al. (2023) and Narechania et al. (2021) allow their participants to ask questions about the task in real time. While these studies have useful insights into how their systems work, they may be more limited in their ability to generalize to average users who might eventually use the system in real contexts.

Tian et al. (2023) and Narechania et al. (2021) also use rule-based template translations directly for simplicity. We observe that SQL queries are easier to translate with templates because they contain more keywords, whereas SPARQL queries emphasize relationships between entities. In SQL, relationships between entities are specified via JOIN operations using foreign keys; spelling these out with templates yields rather unnatural, low-level descriptions.

## 6.   Conclusion

This paper underscores the crucial role of interactivity in semantic parsing for knowledge-based question answering systems, demonstrating through a comparative study and user feedback the heightened effectiveness and satisfaction derived from allowing participants to directly interact with the query translation process. It emphasizes the need for systems to convey complex queries in clear language for user corrections, thus lowering barriers to knowledge base access. The study advocates for human-user interactions and the outlines the potential benefits and drawbacks from crowdsourc-ing.

Future work involves leveraging advanced Large Language Models (LLMs) like ChatGPT, LLaMA, and Gemini to enhance parsing accuracy and user interaction quality (Achiam et al., 2023; Team et al., 2023; Touvron et al., 2023). These models could provide finer interpretations of user queries, thereby enhancing precision. Additionally, developing more intuitive interfaces and feedback mechanisms, coupled with testing systems using real user queries, could offer significant insights into system usability and performance in real-world scenarios. This strategy aims to better meet user needs and understand system efficacy in actual use cases.

## 7.   Ethical Considerations

**Approval from the Institutional Review Board (IRB).** Before starting our crowdsourcing usability study, we secured approval from the Institutional Review Board (IRB) at our institution. This process classified our data collection as Exempt Research, indicating that participation posed minimal or no risk to our human subjects. We did not gather any personal information from participants, except for basic demographic details like their native language, to confirm their eligibility for the task. No identifiable information was collected. Additionally, all participants were required to review and consent to an informed consent form prior to engaging in the study. The identities of MTurk crowdworkers were also anonymized automatically by the platform.

**Paying Crowdworkers.** To ensure quality data collection and fair treatment of crowdworkers, we meticulously formulated our compensation strategy for the MTurk task. Following internal testing, we assessed the average time needed for a task and adjusted the compensation accordingly to meet the minimum wage standards in our state. This adjustment led to a payment of 40 cents per task and $1 per survey. The survey payment was designed to act as both payment and a small bonus. Moreover, to encourage ongoing participation and compensate for the time invested in tutorials and qualification tasks, we provided $5 bonuses upon the completion of each stage in the pipeline.

## 8.   Acknowledgments

# 9. Bibliographical References

Josh Achiam, 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*.

Jacob Andreas and Dan Klein. 2016. Reasoning about pragmatics with neural listeners and speakers. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1173–1182, Austin, Texas. Association for Computational Linguistics.

Shijie Chen, Ziru Chen, Huan Sun, and Yu Su. 2023. Error detection for text-to-SQL semantic parsing. *Findings of the Association for Computational Linguistics: EMNLP 2023*.

Ahmed Elgohary, Saghar Hosseini, and Ahmed Hassan Awadallah. 2020. Speak to your parser: Interactive text-to-SQL with natural language feedback. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2065–2077, Online. Association for Computational Linguistics.

Ahmed Elgohary, Christopher Meek, Matthew Richardson, Adam Fourney, Gonzalo Ramos, and Ahmed Hassan Awadallah. 2021. NL-EDIT: Correcting semantic parse errors through natural language interaction. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5599–5610, Online. Association for Computational Linguistics.

Ben Eyal, Moran Mahabi, Ophir Haroche, Amir Bachar, and Michael Elhadad. 2023. Semantic decomposition of question and SQL for text-to-SQL parsing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13629–13645, Singapore. Association for Computational Linguistics.

Mihir Kale and Abhinav Rastogi. 2020. Template guided text generation for task-oriented dialogue. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6505–6520, Online. Association for Computational Linguistics.

Yunshi Lan and Jing Jiang. 2020. Query graph generation for answering multi-hop complex questions from knowledge bases. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 969–974, Online. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Xintong Li, Symon Stevens-Guille, Aleksandre Maskharashvili, and Michael White. 2021. Self-training for compositional neural NLG in task-oriented dialogue. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 87–102, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Alexander Miller, Will Feng, Dhruv Batra, Antoine Bordes, Adam Fisch, Jiasen Lu, Devi Parikh, and Jason Weston. 2017. ParlAI: A dialog research software platform. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 79–84.

Lingbo Mo, Ashley Lewis, Huan Sun, and Michael White. 2022. Towards transparent interactive semantic parsing via step-by-step correction. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 322–342, Dublin, Ireland. Association for Computational Linguistics.

Will Monroe and Christopher Potts. 2015. Learning in the rational speech acts model. In *Proceedings of 20th Amsterdam Colloquium, Amsterdam.)*.

Arpit Narechania, Adam Fourney, Bongshin Lee, and Gonzalo Ramos. 2021. DIY: Assessing the correctness of natural language to SQL systems. In *26th International Conference on Intelligent User Interfaces*, pages 597–607.

Zheng Ning, Zheng Zhang, Tianyi Sun, Yuan Tian, Tianyi Zhang, and Toby Jia-Jun Li. 2023. An empirical study of model errors and user error discovery and repair strategies in natural language database queries. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*, pages 633–649.

Yilin Niu, Fei Huang, Wei Liu, Jianwei Cui, Bin Wang, and Minlie Huang. 2023. Bridging the gap between synthetic and natural questions via sentence decomposition for semantic parsing. *Transactions of the Association for Computational Linguistics*, 11:367–383.

Sheng Shen, Daniel Fried, Jacob Andreas, and Dan Klein. 2019. Pragmatically informative text generation. 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 4060–4067, Minneapolis, Minnesota. Association for Computational Linguistics.

Elias Stengel-Eskin and Benjamin Van Durme. 2023. Did you mean...? confidence-based trade-offs in semantic parsing. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2621–2629, Singapore. Association for Computational Linguistics.

Alon Talmor and Jonathan Berant. 2018. The web as a knowledge-base for answering complex questions. 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 641–651.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Yuan Tian, Zheng Zhang, Zheng Ning, Toby Li, Jonathan K. Kummerfeld, and Tianyi Zhang. 2023. Interactive text-to-SQL generation via editable step-by-step explanations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16149–16166, Singapore. Association for Computational Linguistics.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Jack Urbanek and Pratik Ringshia. 2023. Mephisto: A framework for portable, reproducible, and iterative crowdsourcing. *arXiv preprint arXiv:2301.05154*.

Ashish Vaswani, 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.

Silei Xu, Shicheng Liu, Theo Culhane, Elizaveta Pertseva, Meng-Hsi Wu, Sina Semnani, and Monica Lam. 2023. Fine-tuned LLMs know more, hallucinate less with few-shot sequence-to-sequence semantic parsing over Wikidata. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5778–5791, Singapore. Association for Computational Linguistics.

Hao Yan, Saurabh Srivastava, Yintao Tai, Sida I Wang, Wen-tau Yih, and Ziyu Yao. 2023. Learning to simulate natural language feedback for interactive semantic parsing. *arXiv preprint arXiv:2305.08195*.

Ziyu Yao, Yu Su, Huan Sun, and Wen-tau Yih. 2019. Model-based interactive semantic parsing: A unified framework and a text-to-SQL case study. 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 5447–5458, Hong Kong, China. Association for Computational Linguistics.

Wen-tau Yih, Matthew Richardson, Chris Meek, Ming-Wei Chang, and Jina Suh. 2016. The value of semantic parse labeling for knowledge base question answering. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 201–206, Berlin, Germany. Association for Computational Linguistics.

## A. Generator Development

### A.1. Motivation

In order for a user to know when a parse is correct, the generated question needs to accurately and clearly reflect the content of the parse, or meaning representation (MR), that it represents. Mo et al. (2022) experimented with various generation models, finding that BART-large performed best and that including dialogue context in the model input improved performance. Furthermore, they find that including templated English versions of the SPARQL queries, which were created to help their crowdworkers understand the meaning of queries in order to rephrase them, also boosted performance.

As mentioned in Section 2.1, we improve the generator by transforming the partial SPARQL query MRs into templated English using a rule-based translation, following the example of Kale and Rastogi (2020). This improvement boosts the generator's performance, though errors still occur. In an analysis of 100 random generated questions, the

| | Specificity (8) | Fluency (1) | Correctness (2) |
|---|---|---|---|
| **Templated** | What is the **composer** of the **music** "My Baby Understands"? | What is the jurisdiction of the governmental body **Hawaii** House of Representatives? | Of which, what is the administrative area of the type UK constituent country? |
| **Human-Written** | Who **wrote** the **song** "My Baby Understands"? | Which government includes the **Hawaii** House of Representatives? | Of the above listed, which one **is a UK constituent** country? |
| **Generated** | Who **wrote** "My Baby Understands"? | What state has the **Hawaraii** House of Representatives as part of its government? | Which of these countries is where **that UK constituent** country is in? |

Table 5: Three instances of errors in an analysis of 100 random generated sub-questions compared to human-written and templated versions. The bold text indicates meaningful differences between them. The numbers in parenthesis indicate the percentage of error cases out of the 100 items examined.

first author finds that errors occur in 11% of them. These errors can be grouped into three categories: specificity, correctness and fluency.

Examples of these three error types can be seen in Table 5. In this table, one can see that under the *specificity* category, the generated question does not mention that the writer is a composer nor that the piece is a song or piece of music, making it less exact in comparison to both the templated question and human-written (gold) question, taken from the INSPIRED dataset. Specificity errors occur in 8 out of 100 instances. *Fluency* errors occur when the generated question has nonsensical grammar or spelling, such as spelling *Hawaii* as "Hawaraii", occurring only once. Lastly, *correctness* errors occur when the generated question does not produce the correct meaning for the LF, which occurs twice.

Though these errors may seem small, they can have a big impact on the usability of the overall system. The specificity error in Table 5, for example, can be interpreted in a number of ways. The question *Who wrote X?* is not specific enough for the user to know that the underlying MR is correct or not, as this could correspond to other closely related meanings such as the writer of a book, lyrics, music, a film, etc. Thus the underlying MR could be incorrect but the user would have a difficult time identifying *how* it is incorrect, in turn making it difficult to provide feedback. The ambiguities and errors that can occur in generated questions are thus important problems to be addressed.



Figure 5: An example of the protocol for experiment 1.

## A.2. Rational Speech Acts (RSA) Modeling

To this end, we employ the linguistic framework of Rational Speech Acts (RSA), which views dialogue as a recursive reasoning process between speakers and listeners to convey meaning (Monroe and Potts, 2015; Andreas and Klein, 2016). The underlying idea is that a listener is reasoning about what "state of the world" is likely to be given that a cooperative speaker produced a particular question, while the speaker is reasoning about how a listener is most likely to interpret the question they produce. Though RSA models focus primarily on the pragmatics of dialogue and thus the information is communicated beyond what is literally said, the concept of recursive reasoning between a speaker and listener can have useful applications for the conveyance of literal meaning in questions by a text generation model.

The work of Shen et al. (2019) demonstrates this effectively by drawing on RSA methods to improve text generations by modeling a listener via *information preservation* (reconstructor-based) and *modeling of distractors* (distractor-based).

## A.3. Experiment 1: Reconstructor-Listener

We adapt the reconstructor-based approach of Shen et al. (2019) to this work, using a base speaker model, $S^0$, to translate MRs to natural language questions and using beam search to output a 10-best list. We then use a reverse model, $L^R$ to represent the listener, which is trained to predict (or reconstruct) MRs from the natural language outputs of the speaker model. We use $L^R$ as a scorer, outputting negative log-likelihood scores for the original MR given each question in a 10-best list from $S^0$: $L^R(i|o)$.

58

| | % of items |
|---|---|
| **Equal Preference** | |
| 1. Listener rank + Speaker rank | |
| 2. Select lowest combined rank (best) | |
| | |
| **Listener Preference** | |
| 1. If Listener rank = 1 and Speaker rank ≤ 3, select | 28.7 |
| 2. Elif Listener rank = 2 and Speaker rank ≤ 3, select | 25.7 |
| 3. Else: | 45.6 |
|     For Listener ranks 1-5: | |
|     Listener rank + Speaker rank | |
|     select lowest combined rank | |
| | |
| **Speaker Preference** | |
| Opposite of Listener Preference | 32.8 |
| | 25.7 |
| | 41.5 |

Table 6: The three methods of selecting a question in each preference model, using rankings. The right side of the table shows the percentages of the items that were selected using that step.

| | Best Option | Breakdown |
|---|---|---|
| Speaker | 17 | C 6 |
| | | S 9 |
| | | F 2 |
| **Listener** | **19** | C 6 |
| | | S 10 |
| | | F 3 |
| Same | 58 | |
| Neither | 6 | |

| Model | Best Option |
|---|---|
| Speaker | 75 |
| Listener | 77 |
| Blend (EP) | 79 |
| Blend (SP) | 80 |
| **Blend (LP)** | **83** |
| Oracle | 94 |

Table 7: (Left) Comparison between pairs of #1 Listener utterances and #1 Speaker utterances of 100 items. Each error was categorized into one of three categories, as seen in the "Breakdown" column. C refers to *correctness*, S refers to *specificity*, and F refers to *fluency*.

Table 8: (Right) Comparison of generated questions selected by different ranking methods on 100 items.

We then use these scores to rerank the questions in order of preference by the listener model. Figure 5 demonstrates this process using an example MR.

In order to assess the differences between the listener model ranking preferences and the speaker model preferences, we select 100 MRs and the corresponding pair of the top ranked speaker and listener questions. We select these items by calculating Levenshtein edit distance between them, filtering out duplicate MRs, and selecting the 100 pairs with the largest edit distance score. The goal of this is to find the cases in which the listener model preference and speaker model preference differ the most, which seems to also signal the cases that are most challenging for the models. We mask the labels of the two questions and randomize them, then manually compare based on differences in correctness, specificity, and fluency (see Table 5 for examples of these criteria). Table 7 shows the results of this analysis. The majority of the examples were of equal quality, though in six instances, neither question reflected the meaning of the original MR. The Listener #1 questions were slightly more successful than the Speaker #1 questions in expressing the meaning of the MR, though the non-negligible success of the Speaker model seems to indicate that there could be benefit to balancing the ranks of the two models.

To this end, we experiment with various methods of blending the scores (as Shen et al. (2019) does) or rankings and find that using rankings has the best results. We experiment with equally weighting the ranks of the two models and preferring one over the other. Table 6 shows the procedure for each of these strategies.

## A.4. Experiment 2: Self-Training

Because the reconstruction-based approach is computationally heavy and unrealistic to use at run time in real dialogues due to latency issues, we experiment with self-training to create a single generation model that is informed by this process.

In the COMPLEXWEBQUESTIONS dataset, there are 24,147 complex questions in the training set that did not get used in the INSPIRED dataset. Thus, we can transform these into our templated representations, giving us more than 50,000 MRs without natural-language counterparts. Using the strategy above, we use the base speaker model to generate a 10-best list of sub-questions for each MR, use the reconstructor model to score them, then use the best blending strategy (listener preference) to select the best single question. We then train a new generation model on this "silver" data for 10 epochs, then 4 epochs on the gold data (from the INSPIRED dataset). We then generate a single sub-question for each of of the 50k MRs using this new model. We iterate this process a second time but find that performance seems to plateau. This is unsurprising; previous work has found that self-training performance usually plateaus after very few iterations (Li et al., 2021). Table 9 shows those results. In this table we use an automatic metric we call *reconstruction accuracy*, which is the percentage of times that the listener model could recover the original MR when given the generated question as input. This gives us a measure of how well the meaning of the MR is preserved in the generated question.

There are a few noteworthy points about these results. First, the speaker, listener, and blend numbers have decreased a bit in this evaluation. The reason for this is basically that, in some cases, even though the question they chose was good, the self-trained model came up with a question that was even better, meaning that it was more specific or context-aware, for example. The self-train number

| Model | Reconstruction Accuracy (%) | Manual Evaluation (%) |
|---|---|---|
| Base Speaker Model | 80.86 | 71 |
| Blended Ranking | 81.74 | 75 |
| Self-Trained | **86.50** | **96** |
| Self-Trained Round 2 | 85.73 | 94 |

Table 9: Comparison of the base speaker model ($S_0$, the best blended rank model (listener preference), and two iterations of the self-trained model.

is also higher than the oracle number in Table 8. This basically means that in two cases, the self-trained model's sub-question was a valid one when the base speaker model did not have a valid question in its list of 10 possibilities. Table 10 shows these two cases.

Because we see substantial gains using the self-trained model at one iteration, we employ this model as our generator in this work.

| #1 utterance | Example 1 | Example 2 |
|---|---|---|
| Template (MR) | What is the **composer** of the music [my baby understands]? | What is the jurisdiction of the governmental body [**Hawaii** House of Representatives]? |
| Speaker | Who **wrote the song** "my baby understands"? | What state has the **Hawaraii** House of Representatives as part of its government? |
| Listener | "My baby understands" is **credited to** whom? | In what state's government can you find the **Hawaraii** House of Representatives? |
| Blend | Who **wrote** "my baby understands"? | In what state's government can you find the **Hawaraii** House of Representatives? |
| Self-Trained | Who is the **composer** of the **song** "my baby understands"? | What state has the **Hawaii** House of Representatives as part of its government? |

Table 10: The two examples where the self-trained model was able to produce a better question than any of the 10-best items produced by the speaker model. Red indicates problem areas, including ambiguity and misspellings, and green indicates contrasting successful areas.

| Question Types | % INSPIRED (N = 3441) | % User Study (N = 400) |
|---|---|---|
| Composition | 43.30 | 48.00 |
| Conjunction | 45.13 | 37.25 |
| Comparative | 6.02 | 8.25 |
| Superlative | 5.55 | 6.5 |
| **Question Facets** | | |
| Filter | 13.02 | 15.00 |
| Restriction | 22.38 | 44.75 |
| Union | 1.60 | 5.00 |
| Sparse Predicates | 10.61 | 15.00 |
| Double Restrictions | 0.44 | 1.00 |
| **Required Edits** | | |
| 0 | 56.81 | 7.50 |
| 1 | 28.74 | 65.25 |
| 2 | 13.28 | 25.75 |
| 3 | 1.08 | 1.50 |
| 4 | 0.09 | 0 |

Table 11: Breakdown of question types, facets, and number of required edits in the INSPIRED dataset and selected questions for the user study (in percentages).

## B. Question Selection

There are 4 main question types: composition, conjunction, comparative, and superlative. There are also many facets possible for a given SPARQL query/question pair, including filters, restriction predicates, and unusual compositions such as union types. We also identify queries that contain a sparse predicate, meaning it appears three or fewer times in the training set of INSPIRED. Note that *predicates* refer to the relations between entities in the knowledge base. See Mo et al. (2022) for explanations of the various question types and facets. Table 11 shows the breakdown of the INSPIRED test set and the 400 selected items within that set used for the user study.

## C. Generation Errors

Figure 6 and 7 show examples of errors made by the generation model. In Figure 6 (which is using the Full system), you can see that the error is quite subtle; the parse correction model predicted the question *Of which, what is the film with **a story by** [Ethan Coen]?* instead of *Of which, what is the film **written by** [Ethan Coen]?* These are in fact different expressions; a film might be based on a story by an author, but the screenplay written by a different author. Because these are very closely related concepts (that are not even made particularly clear by the different SPARQL predicates – *film.film.story_by* versus *film.film.written_by*), the generation model seems to confuse their meanings. This leads to a difference in the underlying SPARQL and template and the natural language

question that the user sees, leading them to believe that the decomposition is correct when it is not.

Figure 7 demonstrates an error by the initial decomposition (in this case using the Decontextualized system) in which the entity name is misspelled, making a sub-question seem incorrect when it is in fact correct.

The first author conducted an error analysis of the dialogues that met the criteria of the *User Deemed Correct* category described in Section 3.1. In those 125 dialogues, 18 were removed due to generation errors. 12 occur in the Full system (across all comparisons, N=96) and 6 occur in the Decontextualized system (N=29). It was considered an error if the natural language question was ambiguous or too vague, for example if it selected the word *location* when the predicate specified *country*, if it contained any misleading language or if it better expressed a different (often closely related) predicate instead of the one it was supposed to be representing (see Figure 6), or if there were fluency or grammatical problems (see Figure 7).

## D. Full vs. Decontextualized System Comparison



Figure 8: A comparison of the success rates for the Full vs. Decontextualized systems (using three different categories – officially correct, inferably correct, and user deemed correct) and the number of edits attempted by the user.

## E. Appendix: User Errors in Templated System

Figure 9 shows a dialogue in the Templated system in which a user achieves the correct decomposition and final answer, but continues to make edits. We assume that this occurs when the user is unsure whether or not a sub-question is correct due to a confusing template.



Figure 6: Generation error by the Full System.



Figure 7: Generation error by the Decontextualized system.

## F. Approximate Randomization Test for Significance

To rigorously test the significance of the observed difference in our metric of effort, we employed approximate randomization testing which is a non-parametric statistical test that does not assume a specific distribution of the data. In this case, this method involves virtually flipping a coin to decide whether to swap the success and edit counts between systems for each dialogue. Through this random reassignment, we generated new sets of data under the null hypothesis that system type does not affect the efficiency metric.

After 10,000 such trials, the proportion of randomized trials yielding a difference in effort scores greater than or equal to the observed difference of 0.05192 was 0.0054 for the Full vs. Templated system. This p-value indicates that the observed difference in effort between the Full and Templated systems is statistically significant, strongly suggest-

**Target Question:**
What is the title of the national anthem of the country where the Vyakta horse breed originated?

Here's how I understood this question:

1. What is the continent that include the country Vyatka horse?
    Ans: NONE

> Replace question #1 with: Which country is the location where the Vyakta Horse breed originated?

Here's what I have according to your feedback:

1. What is the origin of the breed Vyatka horse?
    Ans: Russia

> Insert the following question as question #2: What is the national anthem of this country?

Here's what I have according to your feedback:

1. What is the origin of the breed Vyatka horse?
    Ans: Russia
2. That entity is the country with the national anthem what?
    Ans: National Anthem of Russia

> Change question #2 to: What is the national anthem of that country?

Here's what I have according to your feedback:

1. What is the origin of the breed Vyatka horse?
    Ans: Russia
2. What is the country with the national anthem Vyatka horse?
    Ans: NONE

> ...

Figure 9: Example of an instance in which the template confused a user into thinking a decomposition was incorrect when it was not.

ing that the type of system indeed impacts the efficiency of dialogues.

# Extrinsic evaluation of question generation methods with user journey logs

Elie Antoine[1], Eléonore Besnehard[2], Frédéric Béchet[1,5], Géraldine Damnati[3],
Eric Kergosien[2], Arnaud Laborderie[4]
(1) LIS, Aix Marseille Univ; (2) Geriico, Université de Lille; (3) Orange Innovation;
(4) BNF - Bibliothèque Nationale de France
(5) International Laboratory on Learning Systems (ILLS - IRL CNRS), Montreal

## Abstract

There is often a significant disparity between the performance of Natural Language Processing (NLP) tools as evaluated on benchmark datasets using metrics like ROUGE or BLEU, and the actual user experience encountered when employing these tools in real-world scenarios. This highlights the critical necessity for user-oriented studies aimed at evaluating user experience concerning the effectiveness of developed methodologies. A primary challenge in such "ecological" user studies is their assessment of specific configurations of NLP tools, making replicat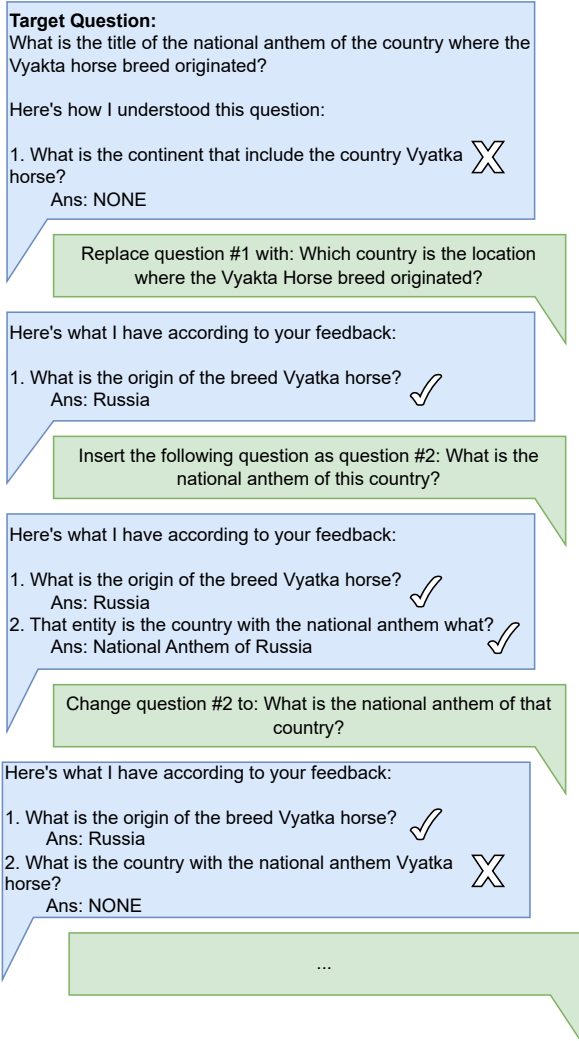ion under identical conditions impractical. Consequently, their utility is limited for the automated evaluation and comparison of different configurations of the same tool. The objective of this study is to conduct an extrinsic evaluation of a question generation system within the context of an external task involving document linking. To do this we conducted an "*ecological*" evaluation of a document linking tool in the context of the exploration of a Social Science archives and from this evaluation, we aim to derive a form of a "*reference corpus*" that can be used offline for the automated comparison of models and quantitative tool assessment. This corpus is available on the following link: https://gitlab.lis-lab.fr/archival-public/autogestion-qa-linking

## 1. Introduction

Question Generation (QG) from text is a key task in Natural Language Processing (NLP), attracting increased attention for its role in testing the syntactic and semantic understanding of generative language models. Recent literature, including Guo et al. (2024), documents the development and comparison of various neural generation techniques, benchmarked against datasets like SQuAD (Rajpurkar et al., 2016) using automated evaluation metrics.

*Intrinsic* evaluations compare machine-generated questions with a human-produced reference set, employing ngram-based metrics such as ROUGE to measure text fluency and semantic similarity metrics like BERTScore, which uses pre-trained BERT embeddings and cosine similarity to assess the closeness of machine and human-generated text.

Additionally, the relevance of human evaluations in assessing question quality is crucial. For example, Bojic et al. (2023) proposes a hierarchical set of criteria for evaluating the semantic content and formulation of Machine Reading Comprehension input questions. As discussed by the authors, these *intrinsic* benchmark evaluations primarily assess question quality, they seldom address the "*usefulness*" of questions in specific applications, indicating a need for *extrinsic* evaluation methods.

We propose to evaluate question generation models through the task of document linking in the general context of exploring Social Science archives with specialized users.

Document Linking consists in adding hyperlinks between documents of a collection according to some criteria. When the criteria are explicit, like in Wikipedia, evaluating the relevance of predicted links means comparing them to a reference containing explicit hyperlinks (Brochier and Béchet, 2021). However when the links are implicit, which is the case when dealing with linking criteria such as *textual similarity* or *entailment*, evaluating links relevance becomes difficult as it relies on subjective criteria and therefore collecting *gold annotation* on such data is a challenge.

In this study, we chose to conduct an experiment using user journey logs [1] to establish sets of related documents during a session. These sets can be used to compare and evaluate different question generation system through the link they produce. Specifically, our aim is to evaluate our question generation system with the "*question-linking*" paradigm as presented in Antoine et al. (2023) with real users by observing the journeys of a panel of testers. These testers explored an interface designed for discovering a collection of journal archives, which offered various exploration options. Among these options, users could select a passage and open a window containing linked passages from other articles.

We compares the links produced by four strategies: the first uses paragraph similarity as a baseline, the second involves similarity between (ques-

---

[1]The corpus collected in this study is available on the following link : https://gitlab.lis-lab.fr/archival-public/autogestion-qa-linking

tion, answer) pairs, with questions generated by a small (< 1B) model and answers extracted from the text. The last two strategies focus solely on question similarity, with one employing a small (< 1B) language model and the other a large (7B) Language Model (LLM).

## 2. Question generation for exploring archive collections

We have explored the potential utility of Question Generation models in the context of exploring a collection of documents. Even if current Question-Answer (QA) models might be too simplistic for use in practical archive exploration, the focus here is on the use of Question Generation models. These models are trained differently from QA models, as they are designed to predict a question based on an answer and a text document, as opposed to generating a response given a question and a document.

The key idea is to use Question Generation models to characterize documents in archives by creating a set of questions associated with the text segments. This is achieved by selecting potential answers from text segments and generating questions based on these answers and their context. By comparing the questions and answers from different documents, the system can predict links between them, effectively adding an *explainability layer* to the document exploration process. This allows users to quickly assess the relevance of links by examining the associated QA pairs, which can save time compared to the traditional approach of following every link to determine its significance.

We present below a short description of our question generation and linking methods.

### Question generation methods
We automatically generated questions on the collection using the same method as the one described in Antoine et al. (2023). In this approach, a semantic parser is used to select potential answers from the articles. As proposed in Pyatkin et al. (2021) and Bechet et al. (2022), a Semantic Role Labelling (SRL) model following the PropBank formalism (Palmer et al., 2005) is used in order to select answers candidates among the detected semantic roles. Following this step, a question generation model is used to provide a question, given the selected answer and its context. This model is trained by fine-tuning the BARThez (Kamal Eddine et al., 2021) language model on a French corpus of question-answer-context triplets called *FQuAD* (d'Hoffschmidt et al., 2020). To address the model's tendency to overgenerate potentially meaningless or overly simplistic questions, a series

of filters are then applied to enhance quality and reduce quantity.These filters are based on resources such as a thesaurus or a list of persons linked to the applicative domains as well as textual indicators. Here is a list of the indicators considered in the filtering process:

1. $\#(pers)$: the number of person mentions belonging to a given list

2. $\#(th\_answer)$ and $\#(th\_question)$: the number of keywords from the thesaurus of notions in respectively the answer and the question (to add a control on the semantic relevance of the question)

3. We compute the average length of the generated questions and extracted answers to calculate the deviation from the mean of each questions ($quest\_diff\_mean$) and answers ($ans\_diff\_mean$)

4. We finally compute $inter\_qa$, the percent of intersection between the extracted answer and the question (to avoid nonsensical questions that contain the answer to their own question).

All these filters are used in a decision rule that accept or reject a generated pair question/answer.

### Linking methods
Links between items in the collection are produced using the same method as in Antoine et al. (2023). The proposed approach is to generate links using questions and answers generated from the text rather than directly on the text itself. The embedding projection for each "`<question> | <answer>`" pair structure uses the Sentence-Transformer (Reimers and Gurevych, 2019) library, and more precisely the multilingual model *distiluse-base-multilingual-cased-v1* (Reimers and Gurevych, 2020). A cosine similarity measure is then employed between all pairwise combinations of these embeddings, resulting in the computation of a similarity matrix.

In this study we will perform an extrinsic human evaluation where the usefulness of the questions for document linking is studied.

## 3. Collecting logs from an exploration interface

This study was conducted within the framework of the French ANR project ARCHIVAL[2], aimed at developing novel exploration methods for thematic archive collections using machine comprehension techniques. The archive collection chosen for this study is a collection of social science journal articles

---

[2] https://anr.fr/Projet-ANR-19-CE38-0011

in French from the *Autogestion* (Self-management) journal[3]. This collection is distributed in its digitized form by the French Persée organization. It is part of a larger pluridisciplinary multilingual mixed collection (archives and documents) that has been gathered since the 1960's by the FMSH[4] foundation's library. The full collection has been granted the Collex label (*Collection d'Excellence* or Excellency Collection) from the CollEx-Persée[5] network under the supervision of higher education and research for the preservation of corpus of digitized or natively digital documents.

This collection, published during a period ranging from the 1960s to the 1980s, constitute a corpus of 46 issues for an overall amount of 896 articles (more than 6000 pages and 1.98M tokens).

In order to navigate in this large collection, the interface homepage proposes a search engine and two main access modes: a direct access through timelines, tables of content and indexes containing references to persons (all the authors and people mentioned in the documents content) and notions from a thesaurus of around 400 notions specifically designed for the semantic domain of the journal. The notions and persons are automatically detected from the text of the articles, the method and the exploitation of these functionalities by users are out of the scope of the current study and are not detailed here.

Once a user has entered the collection and opened an article, he can further explore it with linking mechanisms. The user can select a text area in the article which becomes highlighted. This selection corresponds to a particular area of interest for which links to other documents in the collection can be proposed according to two methodologies:

- Firstly question-linking method presented in the previous section. A list of questions generated from the paragraph containing the highlighted text is displayed to the user who can click on any of these questions to obtain a list of $n$ links to related paragraphs in other documents calculated thanks to the method presented before. Links are associated here to references to the title and the authors of the target documents as well as a snippet of the target paragraph. The link is explained by the pair of questions from the source and the target paragraph. An example of document linking and question explanation is given in figure 1.

- Secondly, a method based on textual similarity using *SentenceBert* (Reimers and Gurevych,

2019) is applied to the paragraph containing the highlighted text in order to display the $n$ other paragraphs in the collection that minimize the similarity criteria.

In our experiments, the amount of displayed links was set to $n = 10$. An illustration of this text selection and linking presentation method is given in figure 2. Users can choose the document linking method they want to use.

Within the interface, users can perform a variety of actions to navigate and manipulate content. First, they can open or close windows associated with articles, notions, or persons.They can also switch between different views, including the timeline, notions page, and persons page.They can switch submenus within an article window, whether it's toggling between viewing the article text, the notions automatically extracted from the text or the person cited in the text.

To seek relevant connections, users can also use the links provided when selecting an area of interest. All these actions are logged.

A first way of exploiting the logs would be to analyse if users actually clicked on the links proposed by the various algorithms. If this is an interesting way to analyze user journeys and their acceptation of the functionalities, it is not enough to provide a reproducible evaluation framework to compare several question generation approaches or several linking strategies. In this work, in order to propose a reproducible evaluation protocole, we consider that the set of documents consulted by a given user during a test session constitutes a coherent set of documents that are of interest for him/her. We will call this set of consulted documents a *user-log collection*. Then we want to check a posteriori if, starting from one document of the collection, a given exploration approach would allow to reach other documents from the same collection. We formulate the hypothesis that proposing links that allow users to reach more easily other documents of interest is more helpful. Hence we can compare several linking methods, beyond the ones that were originally implemented during the collection phase.

## 4. From log collection to extrinsic evaluation

This section describes how we turn the set of documents in our corpus into a graph according to a given linking method, and how we can evaluate such graphs thanks to the user-log collections described earlier.

**Graph creation**

For each linking method $L$, the first step in our process is to turn our document collection into a

---

Figure 1: Example of highligthed source and target paragraph with question explanation



Figure 2: Text selection and link production interface

graph $G_L$. This is achieved by considering all the documents within our corpus which contain automatically generated links. Each node in $G_L$ corresponds to a document (an article of the *Autogestion* journal), and we add an edge between document $A$ and document $B$, noted $(A, B)$, if there is at least one link connecting a paragraph from $A$ to a paragraph in $B$ thanks to the linking method $L$. This is a directed graph as all linking methods are not necessarily symmetrical.

We apply a weight to all the $(A, B)$ edges of this graph between a document $A$ and a document $B$ according to the following principle:

1. for each edge $(A, B)$ we compute the number of direct links between documents $A$ and $B$, called $N_L(A, B)$

2. to normalize these numbers at the document level, for each document $A$, we rank all the outgoing edges from $A$ to any other document $(A, .)$ in the collection according to the values $N_L(A, .)$.

3. the weight of edge $(A, B)$ called $W_L(A, B)$ is the rank of this edge among all the outgoing edges from document $A$ sorted by $N_L(A, .)$.

The best weight an edge $(A, B)$ can have is $W_L(A, B) = 1$, corresponding to the pair of documents having the highest number of links according to the linking method $L$. The worst weight for $W_L(A, B)$ is the maximum number of outgoing edges from $A$ (bounded by the number of documents in the collection).

The document graphs obtained for each linking

Figure 3: Example of number of jumps (without weight for clarity) between documents in a user log, compared between two linking methods on the collection. Note that because our graph is directed, the number of jumps between two documents, e.g. A and B (5 for method 1) will be different than between B and A (6 for method 1).

method are evaluated by their ability of visiting efficiently, by taking into account the weights previously described, each user-log collection of documents. As we can see in figure 3 where no weight is specified for clarification purpose, to visit all documents selected by user 1, we need to follow on average $4.83$ links with the linking method 1 and $2.75$ with method 2. Therefore, we will consider that the linking method 2 is more efficient for recreating the journey of user 1 than method 2.

**Metric**

To evaluate the quality of our different methods on the graph produced, we use the average number of jumps corresponding to the average length in terms of edges in the weighted shortest paths between all node pairs in a set of documents.

This metric can be viewed as the number of links to be clicked on in a source document to reach the target article, or the number of intermediate articles to be visited, counting the initial one, as shown in figure 3. Since weights in our graph are in increasing order of importance (the best weigh is 1), finding the length of the weighted shortest path between documents $A$ and $B$ gives insight into the likelihood of a reader navigating from $A$ to $B$ via recommendation links.

We can then compare these values across methods and juxtapose them with the average number of weighted jumps between all document pairs in the collection.

## 5. Experiments

### 5.1. User log collection

Three test sessions were held in May, June and November 2023 to test the *ARCHIVAL* demonstrator with potential users. During the first test day a total of sixteen testers came together for a experiment of discovery and familiarization with the *ARCHIVAL* system. Two panels were set up to carry out two test sessions: the morning session brought together ten testers, mainly researchers in information and communication sciences, while the afternoon session was made up of six testers with a profile of library and documentation professionals. Four testers were invited for the second session of experiment: two information and communication sciences teacher-researchers, a PhD student and teacher-documentalist, and a librarian. Then the third session gathered five expert researcher in the domain of the OCRized journal. For each test session, the general framework was the same and we consider all testers to be part of a single panel of 25 users.

During each test session, testers were instructed to explore the demonstrator freely. After 40 minutes of free exploration where we observed their use of the interface, they answered an initial general questionnaire on their apprehension and appropriation of the device, their use of certain functionalities and their documentation habits. The testers then continued the experiment using suggested entry articles.

## 5.2. Linking strategies

We performed question generation and linking on the 896 documents used in this study. The average number of questions generated by the BARThez model for each granularity level are given in table 1 as well as the percentage of elements containing at least one question for each level. We can see that about 16% of the documents do not contain any question, this corresponds mainly to the summaries or bibliography where we could not generate questions. Less than half of all paragraphs contain at least one question, with an average of 1.0 questions per paragraph and 2.7 if we exclude paragraph with no questions at all. The 60.4% of paragraph that doesn't contain any question consists either of very short ones such as end notes, titles and all micro-textblocks detected by the OCR or of paragraphs where our question filtering process discarded all the questions generated as being non relevant.

For comparison, a second question generation method based on a larger model, Mistral-7B-Instruct-v0.2 (Jiang et al., 2023) was employed. This time, generation was conducted directly at the page level, without prior extraction of potential answers, utilizing an empirically created prompt (as see in listing 1) and no additional filters. Pages were selected based on the presence of at least one question generated by the BARThez model for generation with Mistral-7B-Instruct-v0.2. We aligned the number of generated questions to the one obtained with the previous approach, a portion of questions were randomly removed to ensure balanced comparison. Subsequently, an average of 71.6 questions were generated per article using this method.

Listing 1: Mistral prompt

```
You're a professor of history in the
    field of human and social
    sciences. Annotate the document
    in the form of open questions in
    French as you read about key
    elements of the given paragraph.
    The questions shouldn't be too
    verbose, and may relate to
    elements whose answers are
    present in the paragraph or not.
{document}
Questions :
-
```

Following the methodology in section 5 We have generated four document linking graphs, as described in figure 3, one for the question-linking method $G_{qa}$ using both questions and answers to compute similarity measures, and one for the paragraph similarity method $G_{para}$. The two other ones correspond to the graphs produced by the same method, applied only on the questions of BARThez

| Measure | Article | Page | Paragraph |
|---|---|---|---|
| avg. nb. Q. per element | 70.4 | 9.4 | 2.7 |
| % elements with Q. | 83.8% | 84.6% | 39.6% |

Table 1: Average number of questions generated at each level of granularity (document, page, and paragraph) for the BARThez model and percentage of articles with at least one question



Figure 4: Count of document by number of incoming and outgoing links for $G_{qa}$

$G_q$ and Mistral $G_{mistral}$, the latter having no answer extraction.

To build $G_{qa}$ we computed the cosine similarity metric between all the SentenceBert (Reimers and Gurevych, 2019) embeddings of the concatenation of the question and answer (question + answer) and kept the top 10 links for each of them. The question where generated on 896 articles (corresponding to nodes) and produced 136,478 unique links in total (corresponding to edges).

For $G_{para}$ we computed the cosine similarity metric using the embedding of the paragraphs having produced at least one question. We generated the links with the same constraints as for the questions + answer, with a maximum of 10 links. We produced 129,550 unique links for this methods.

Finally, for $G_q$ and $G_{mistral}$ we computed the cosine similarity metric directly on the embeddings of the generated questions of both models. We generated the links with the same constraints as for the questions + answer, with a maximum of 10 links. We produced respectively 147,662 unique links for $G_q$ and 153,454 unique links for $G_{mistral}$.

## 5.3. Results

All our constructed graphs feature a single, strongly connected component. This result shows that it's possible to explore the entire collection using the links produced by all the used methods, without getting stuck in a clique.

In our experiments we have 25 users, so we used 25 user-log collections, with an average of 13.7 documents in each set.

| | #links | Average number of jumps | |
|---|---|---|---|
| Graph | avg. | All articles | User logs |
| $G_{para}$ | 144 | 5.18 | 4.34 |
| $G_{qa}$ | 152 | 4.88 | 4.01 |
| $G_q$ | 165 | 4.69 | 3.83 |
| $G_{mistral}$ | 174 | 5.26 | 4.34 |

Table 2: Average number of links (in and out) for each method and average number of jumps for all pairs of articles in the entire collection (All articles) and in user-log collections

The average number of unique links in and out of each document is given in table 2. A more precise breakdown of articles according to their number of incoming and outgoing links for $G_{qa}$ is shown in figure 4.

We can see in table 2 that for all methods, the average number of jumps between the users articles is lower than the average number of jumps between articles in the collection. We can assume that those methods gives the user easier access to articles considered relevant than to a random article, the link using question being the one bringing explored articles closer together.

BARThez's question-only linking method gives the best results over the other results, and specifically over his question+answer variant, with the shortest average path. This result is consistent with feedback from platform users who told us that they didn't find the answer useful in their search for links, and that it could even confuse them.

The questions produced by Mistral do not yield links as dense as the other question generation methods, with scores close to the one of the similarity between paragraph. This can be explained by several factors, the first being the granularity of the generation, at page level rather than paragraph level. The second is the generation method and task, with prompting for more open-ended and general questions than SQuAD-style text comprehension questions with already-defined answers. The last is the absence of an expert filter on question generation, as described in section 2.

These experiments show that it is possible to use logs from users exploration in order to compare and evaluate linking methods as an extrinsic task for evaluating the usefulness of question generation methods. The results obtained can give some indications about the efficiency of finding connected documents with a given linking method.

## 6. Conclusion

In this paper, we introduced a framework and approach for harvesting ecological user logs to evaluate a question generation method trough an extrinsic document linking task. By exploiting graph metrics, we conducted evaluations using these logs to gain insights into the links generated by our method, comparing them with links produced by a LLM and traditional linking techniques. Our results highlight a notable observation: even a compact model such as BARThez, enhanced with expert filters and heuristics, can outperform a generic-purpose LLM in generating task-specific questions. This underscores the effectiveness and robustness of our methodology in enabling a comparison of questions through an extrinsic document linking task, offering insights into the efficacy of various question generation approaches trough this specific task. The data collected in this study is available on the following link: https://gitlab.lis-lab.fr/archival-public/autogestion-qa-linking.

Elie Antoine, Hyun Jung Kang, Ismaël Rousseau, Ghislaine Azémard, Frédéric Bechet, and Géraldine Damnati. 2023. Exploring social sciences archives with explainable document linkage through question generation. In *Proceedings of the 7th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*, pages 141–151.

Frederic Bechet, Elie Antoine, Jérémy Auguste, and Géraldine Damnati. 2022. Question generation and answering for exploring digital humanities collections. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4561–4568, Marseille, France. European Language Resources Association.

Iva Bojic, Jessica Chen, Si Yuan Chang, Qi Chwen Ong, Shafiq Joty, and Josip Car. 2023. Hierarchical evaluation framework: Best practices for human evaluation. *Human Evaluation of NLP Systems*, page 11.

Robin Brochier and Frédéric Béchet. 2021. Predicting links on wikipedia with anchor text information. In *Proceedings of the 44th International*

*ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '21, page 1758–1762, New York, NY, USA. Association for Computing Machinery.

Martin d'Hoffschmidt, Wacim Belblidia, Quentin Heinrich, Tom Brendlé, and Maxime Vidal. 2020. FQuAD: French question answering dataset. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1193–1208, Online. Association for Computational Linguistics.

Shasha Guo, Lizi Liao, Cuiping Li, and Tat-Seng Chua. 2024. A survey on neural question generation: Methods, applications, and prospects. *arXiv preprint arXiv:2402.18267*.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.

Moussa Kamal Eddine, Antoine Tixier, and Michalis Vazirgiannis. 2021. BARThez: a skilled pre-trained French sequence-to-sequence model. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9369–9390, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. The proposition bank: An annotated corpus of semantic roles. *Comput. Linguist.*, 31(1):71–106.

Valentina Pyatkin, Paul Roit, Julian Michael, Yoav Goldberg, Reut Tsarfaty, and Ido Dagan. 2021. Asking it all: Generating contextualized questions for any semantic role. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1429–1441, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392. Association for Computational Linguistics.

Nils Reimers 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.

Nils Reimers 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 (EMNLP)*, pages 4512–4525, Online. Association for Computational Linguistics.

# Towards Holistic Human Evaluation of Automatic Text Simplification

**Luisa Carrer[1], Andreas Säuberli[2], Martin Kappus[1], Sarah Ebling[2]**
[1]School of Applied Linguistics, ZHAW Zurich University of Applied Sciences
[2]Department of Computational Linguistics, University of Zurich
cars@zhaw.ch, andreas@cl.uzh.ch, kapm@zhaw.ch, ebling@cl.uzh.ch

## Abstract

Text simplification refers to the process of rewording within a single language, moving from a standard form into an easy-to-understand one. Easy Language and Plain Language are two examples of simplified varieties aimed at improving readability and understanding for a wide-ranging audience. Human evaluation of automatic text simplification is usually done by employing experts or crowdworkers to rate the generated texts. However, this approach does not include the target readers of simplified texts and does not reflect actual comprehensibility. In this paper, we explore different ways of measuring the quality of automatically simplified texts. We conducted a multi-faceted evaluation study involving end users, post-editors, and Easy Language experts and applied a variety of qualitative and quantitative methods. We found differences in the perception and actual comprehension of the texts by different user groups. In addition, qualitative surveys and behavioral observations proved to be essential in interpreting the results.

**Keywords:** automatic text simplification, Easy Language, post-editing, human evaluation, reading comprehension

## 1. Introduction

Text simplification is a form of intralingual translation, or rewording, within one language, i.e. from the standard variety into a simplified variety (cf. Hansen-Schirra et al., 2020). Easy Language and Plain Language are two examples of easy-to-understand varieties aimed at optimizing the readability and comprehensibility of texts for a wide and heterogeneous target audience. More specifically, Easy Language is a strongly controlled form of language and is based on strict sets of rules (Maaß, 2020; Bock and Pappert, 2023). Primary target groups include persons with intellectual disabilities, persons with functional illiteracy, L2 learners as well as persons with dementia, prelingual hearing impairments, and aphasia (Bredel and Maaß, 2016). As a natural language processing task, automatic text simplification (ATS) has increasingly gained traction in recent years (Štajner, 2021). However, there is no consensus on best practices for evaluating simplified texts, resulting in inconsistencies in the methods used (Grabar and Saggion, 2022). Most commonly, automatic evaluation metrics are used, which have been shown to be unreliable (Alva-Manchego et al., 2021).

Studies that involve human evaluation typically employ experts or crowdworkers to rate different aspects of the output text such as simplicity, fluency, and adequacy on Likert-style scales (Štajner, 2021). However, those approaches have several shortcomings: first, they are not representative of the primary target groups of simplified texts. Second, they do not include other stakeholders such as post-editors. Third, they heavily rely on sub-

jective ratings, which may not be indicative of the functionality of the simplified texts, i.e., enhanced comprehensibility.

In this paper, we contribute to the current debate on best practices for human evaluation by exploring different ways of measuring the quality of automatically simplified texts. Our methods span the quantitative to the qualitative, the subjective to the objective, and our raters range from Easy Language professionals to end users. Specifically, we conduct three evaluation studies: an end-user comprehensibility evaluation (Section 3), a post-editing productivity study (Section 4), and an expert evaluation (Section 5). Finally, we discuss the benefits of such multi-faceted evaluations of ATS and provide recommendations for future work.

## 2. Background and Related Work

### 2.1. Human Evaluation of Text Simplification

In terms of human evaluation, previous research has primarily relied on Likert-scale ratings of simplicity, fluency, and adequacy or meaning preservation for evaluating the quality of ATS output (Al-Thanyyan and Azmi, 2021; Stodden, 2021; Ryan et al., 2023; Martin et al., 2022; Štajner and Nisioi, 2018; Mallinson et al., 2020). The raters in these studies are typically researchers, students, or crowdworkers.

Štajner (2021) argued that evaluating ATS output quality should include the usability by target readers. However, evaluation including target groups of Easy Language are rare. Notable ex-

ceptions include studies involving deaf and hard-of-hearing adults (Alonzo et al., 2021), persons with intellectual disabilities (Huenerfauth et al., 2009; Saggion et al., 2015) or dyslexia (Rello et al., 2013b,a,c), and language learners (Crossley et al., 2014). In some cases, comprehensibility is assessed based on comprehension tests, e.g., using multiple-choice questions (Leroy et al., 2013, 2022; Fajardo et al., 2014; Charzyńska and Dębowski, 2015; Alonzo et al., 2021), cloze tests (Charzyńska and Dębowski, 2015; Redmiles et al., 2019), or free recall questions (Leroy et al., 2013, 2022). More rarely, measurements of reading behavior such as reading speed (Alonzo et al., 2021; Crossley et al., 2014; Saggion et al., 2015; Rello et al., 2013a), scrolling interactions (Gooding et al., 2021), or eye movements (Rello et al., 2013a,c) are obtained.

## 2.2. Evaluation of Post-editing Effort

The widespread use of post-editing in interlingual translation has spurred significant research interest in how translators engage in this task and the level of effort involved. Since Krings' (2001) seminal work, it has been widely recognized that post-editing effort encompasses three main dimensions: temporal, technical, and cognitive (cf. Alvarez-Vidal and Oliver, 2023). Temporal effort is easily quantifiable and directly influences productivity and is thus used to determine translators' post-editing rates. Technical effort pertains to the editing actions performed during post-editing, such as text productions, text eliminations, replacements, and shifts, often analyzed using keylogging data and specialized software. Finally, cognitive effort refers to the mental processes underlying post-editing, even when no tangible changes are made to the raw machine translation (MT) output. Measuring cognitive effort is challenging due to its complexity, but pauses have emerged as indicative of cognitive load. Lacruz et al. (2012, 2014) proposed measuring clusters of short pauses, which revealed a clear correlation with post-editing effort. To the best of our knowledge, the present paper represents the first evaluation of post-editing effort for text simplification.

# 3. End-user Evaluation

In this section, we describe an evaluation involving two groups of end users (with and without intellectual disabilities). We measured text comprehensibility with comprehension questions and perceived difficulty of automatically simplified German texts and compared those measurements to the original (non-simplified) source texts and manually created reference simplifications of those texts. The end-user evaluation was already described in more detail in Säuberli et al. (2024) and will only be summarized here.

## 3.1. Materials and Methods

### 3.1.1. Texts and Comprehension Questions

The texts we used in this study are part of a parallel corpus of original and simplified German texts. The corpus was made available to us by a commercial provider of text simplification services in the context of a large-scale research project on automatic text simplification. The texts span various topics and genres, including news, administrative texts and political advertisements. Their lengths range between 100 and 600 words.

Each text exists in three versions: (1) the original source text, (2) a reference simplification, which was manually created by the provider, and (3) an automatically simplified version. We generated the latter with a transformer-based model fine-tuned on data from the same parallel corpus using the approach described in Rios et al. (2021).

Based on the source and reference texts, we created four multiple-choice comprehension questions for each of the 12 texts. One of the questions was about the overarching topic of the text (with four answer options), while the remaining three asked about specific details in the text (with three answer options each).

Since the ATS model sometimes omits information from the source text, and the comprehension questions were written only based on the source and reference texts, some of the questions are not answerable based on the automatically simplified version. Therefore, we added a fourth answer option "Information does not appear in the text" to the detail questions.

### 3.1.2. Participants

To compare comprehensibility among different populations, we recruited two groups of participants. The target group consisted of 18 persons with intellectual disabilities, i.e. a primary target group of Easy Language. The control group consisted of 18 native German speakers without intellectual disabilities. All participants took part on a voluntary basis and were compensated monetarily.

### 3.1.3. Procedure

Data collection was conducted using a mobile app which allowed participants to read and rate the texts and answer the comprehension questions. The texts were randomly assigned to participants such that each participant read exactly one version of each of the 12 texts.

After reading a text, participants were asked to rate the difficulty of the text on a five-point scale

(a) Comprehension question responses



(b) Perceived difficulty ratings

Figure 1: Difficulty estimates of the text versions based on responses to the comprehension questions and ratings. Points are posterior medians, error bars are 80%, 90% and 95% credible intervals (CI). A bracket with ▲ indicates that the 80% CI of the difference between the two parameters does not include zero (i.e., we are 80% confident that there is a difference). Similarly with ▲▲ for 90% CI and ▲▲▲ for 95% CI.

(from "very difficult" to "very easy"). The text was then shown again, and the comprehension questions were displayed at the same time. Apart from the responses to the rating and comprehension questions, we also collected behavioral data, including reading speed and response time.

### 3.2. Results

In this section, we will highlight the results for comprehension questions and difficulty ratings. Refer to Säuberli et al. (2024) for more detailed results.

To estimate the difference in comprehensibility between the three text versions based on the responses to the comprehension questions and the ratings, we applied Bayesian one-parameter logistic item response models (also known as Rasch models; Fox, 2010) and modeled the difficulty of the text version as an additional latent trait (cf. Linacre, 1989). We used separate models for each group.

The difficulty estimates based on the comprehension questions and the difficulty ratings are shown in Figure 1. For the target group, the ATS output was not significantly different from the source or the reference in terms of difficulty. For the control group, the output was even slightly more difficult than the source. The difficulty ratings also show remarkable differences between the two groups. While the target group rated the ATS output as being simpler than the source text, the control group's ratings suggest that the output was equally difficult

as the source and significantly more difficult than the reference.

### 3.3. Discussion

Several remarkable differences can be observed between the results of the two groups. First, the estimated effects are smaller and more uncertain in the target group. This is likely due to the heterogeneity of the target group, but also due to noisier data. The behavioral measurements show that reading speeds varied widely within the target group, suggesting that some participants did not read the texts carefully before rating, leading to less predictable responses. This interpretation is supported by the fact that the target group's difficulty ratings did not differ significantly between the source and reference texts, while the control group's ratings did.

## 4. Post-editing Evaluation

### 4.1. Methods

In our post-editing (PE) evaluation, we observed human translators post-editing the output of the automatic text simplification model (cf. Section 3.1.1) and we quantified PE effort. The following methods were employed: (a) a pre-task questionnaire to collect professional background data as well as attitudinal data on participants' practices in both inter-

and intralingual translation settings; (b) automatic recording of participants' unfolding typing process in manual simplification (MS) vs. PE tasks; (c) a post-task questionnaire to investigate how the participants rated their productivity during the tasks; and, finally, (d) a comparative analysis of production time and effort required.

In line with common practice in translation process (cf. Alves, 2003; Kappus and Ehrensberger-Dow, 2020) and post-editing research (cf. Krings, 2001; Alvarez-Vidal and Oliver, 2023), several quantitative measures were used to determine the effort involved in manually simplifying the eight source texts and in post-editing the corresponding automatically simplified target texts. More specifically, effort was quantified in terms of task duration, number of keyboard and mouse-based user events (as a measure of addition, change, regression or navigation), number of cognitive pauses (i.e., pauses with duration greater than 2000 milliseconds), and total pause time.

## 4.2. Participants and Procedure

Four German-speaking professional translators from a commercial provider of text simplification services were recruited through self-selection sampling. Each participant was given detailed step-by-step instructions to perform two MS and two PE tasks in their workplace. Eight source texts were used in this phase of the study. Texts were selected from an *ad-hoc* pool of texts used for all human evaluations (cf. Section 3.1.1). To prevent bias during the PE activity, each participant manually simplified and post-edited two different pairs of texts. Keystroke logging (GenoGraphiX-Log; cf. Caporossi et al., 2023) and screen recording of both processes were employed.

## 4.3. Results

Data from the pre-task questionnaire showed that participants (P) had three to five years of professional experience in text simplification and various degrees of expertise in interlingual MT and PE, with responses spreading evenly across choices (from 'no experience' to '3-5 years' of experience). On the other hand, participants' background in automatic intralingual text simplification and PE was significantly lower, with three out of four respondents having less than one year of experience.

Table 1 shows the total and mean values for each measure of effort compared between the MS and PE tasks. Student's t-tests with unequal variances were used for statistical analysis. All statistical tests were one-tailed with a 5% level of significance (p < 0.05).

As can be seen from Table 1, no significant difference could be determined for any of the effort measures considered. This means that the statistical data do not suggest any significant decrease in effort in either MS or PE activity.

However, it should be emphasized that borderline statistical values relating to three user events categories were extracted (see Table 2 for a descriptive user events analysis). The mean number of text productions (i.e., textual inputs) in the MS tasks was higher (i.e., 2105, range of 952–3429) than the mean number of text productions in the PE tasks (i.e., 1417, range of 513–3043; t = 1.77; p = 0.06). On the other hand, the mean number of cursor navigations (i.e., navigation key presses) in the MS tasks was lower (i.e. 118, range of 4–548) than the mean number of cursor navigations in the PE tasks (i.e. 724, range of 4–2265; t = 1.85; p = 0.05). Similarly, study participants in the MS tasks made on average fewer mouse clicks (i.e., 55, range of 2–92) than they did in the PE tasks (i.e., 94, range of 7–226; t = 1.81; p = 0.09).

In the post-task questionnaire, respondents were asked to rate their perceived productivity on a 5-point Likert scale (5 = very high). Self-assessed productivity reached an average rating of 4.25 in the MS tasks. In their comments, respondents reported that they could generally maintain a high concentration during the task and that the source texts were "readily comprehensible". In the PE tasks, self-assessed productivity reached an average rating of 3.00. Three out of four respondents were unanimous in pointing out that "cognitive pauses" were often necessary, ultimately affecting productivity. Despite admitting that the automatic output provided a helpful "rough structure" and seemingly good translation solutions, respondents reported that the target texts lacked coherence. In addition, they emphasized that a painstaking source-target comparison was necessary to validate the adequacy of the automatic output, which resulted in higher time expenditure and lower productivity. On the other hand, one respondent stated that the post-editing activity required little effort overall, as the source texts were "relatively easy".

## 4.4. Discussion

The mean productivity value that participants self-reported was 1.25 points higher in the MS tasks than in the PE tasks. Nonetheless, our statistical analysis did not suggest any significant increase in productivity in either manual simplification or post-editing activity. Factors that may have affected the results include participants' main expertise in manual simplification (vs. post-editing). Furthermore, a different working environment – that is, the use of keystroke logging software to perform the MS and PE tasks – may have had an impact on the participants' translation activity and/or their perceived productivity.

| Effort measure | Mean (per P per task) | | Student's t-test |
|---|---|---|---|
| | MS | PE | |
| Task duration (h:m:s) | 00:21:15 | 00:22:14 | p = 0.43 \| t = 1.78 |
| User events | 2760* | 2556 | p = 0.39 \| t = 1.79 |
| Number of cognitive pauses | 87.57* | 88.25 | p = 0.48 \| t = 1.78 |
| Pause time (h:m:s) | 00:18:47* | 00:19:22 | p = 0.45 \| t = 1.77 |

Table 1: Total and mean values for each measure of effort compared between manual simplification (MS) and post-editing (PE) tasks (* P1 completed both MS tasks but did not submit keystroke logging data for the second MS task). Right column: statistical analysis of each measure of effort in manual simplification (MS) and post-editing (PE) tasks according to two-sample Student's t-tests assuming unequal variances.

| User events | P1 | | P2 | | P3 | | P4 | |
|---|---|---|---|---|---|---|---|---|
| | MS* | PE | MS | PE | MS | PE | MS | PE |
| Text productions | 1708 | 1092 | 5697 | 3064 | 2772 | 1308 | 4559 | 4665 |
| Text eliminations | 87 | 428 | 894 | 767 | 321 | 122 | 773 | 1450 |
| Cut events | 0 | 1 | 3 | 7 | 2 | 1 | 2 | 3 |
| Copy events | 0 | 2 | 1 | 1 | 0 | 3 | 1 | 0 |
| Paste events | 0 | 4 | 6 | 8 | 2 | 6 | 12 | 4 |
| Cursor navigations | 153 | 3171 | 16 | 66 | 95 | 9 | 567 | 2553 |
| Mouse events | 2 | 23 | 156 | 274 | 80 | 124 | 150 | 333 |
| Misc. events (e.g., modifier keys) | 133 | 136 | 567 | 336 | 181 | 83 | 383 | 407 |
| Total user events | 2083 | 4857 | 7340 | 4523 | 3453 | 1656 | 6447 | 9415 |

Table 2: Events analysis per study participant (P) in manual simplification (MS) and post-editing (PE) tasks. Combined values (i.e., two MS tasks and two PE tasks) per study participant (* P1 completed both MS tasks but did not submit keystroke logging data for the second MS task).

## 5. Expert Evaluation

### 5.1. Methods

In our expert evaluation, we obtained translation quality ratings from experts in German Easy Language translation. In this phase, we employed an online evaluation questionnaire in which four evaluators performed a source-based direct assessment (cf. Graham et al., 2013; Federmann, 2018) of the target texts. The questionnaire was developed with LimeSurvey[1] and comprised eleven items, of which three collected professional background data, and eight presented two parallel texts each, i.e. one source text and one corresponding target text. For each source text used in the post-editing productivity study (cf. Section 4.2), four corresponding simplified versions were employed, i.e. one reference text, one automatically simplified text, one manually simplified text, and one post-edited text – the latter two being produced during the post-editing study (cf. Section 4). Based on the experts' evaluations, the end quality of the experimental units was then analyzed and compared.

### 5.2. Participants and Procedure

Four Swiss-based German-speaking experts in Easy Language translation were recruited through purposive sampling. A 4x8 Latin square gave us a total of 32 experimental units and secured an unbiased response. Evaluators were asked to assign simplicity, adequacy, and fluency scores on 5-point scales (5 = maximal quality; cf. Grabar and Saggion, 2022) to each target text (see Table 3) and, if desired, insert comments. Evaluators were not provided with any information about how the target texts had been produced.

### 5.3. Results

All four evaluators had over five years of professional experience in text simplification and regularly provided a wide portfolio of Easy and Plain Language services, including intra- and interlingual translations, text production, and text validation in collaboration with the target groups. A comparison of the average simplicity ratings for each target text category shows that the manually simplified texts produced during the post-editing productivity study (cf. Section 4) were rated higher (i.e., 4.38, range of 4–5) than the other three categories (see Table 4). A similar pattern emerged for the adequacy and fluency ratings: the manually simplified texts

---

[1] https://www.limesurvey.org/

| Rating | Simplicity Q1: How does the target text differ from the original text? | Adequacy Q2: Does the target text reflect the content of the original text? | Fluency Q3: Is the target text fluent and grammatical? |
|---|---|---|---|
| 5 | much easier | completely | fluent |
| 4 | easier | mostly | mostly |
| 3 | equally difficult | partially | partially |
| 2 | more difficult | mostly not | mostly not |
| 1 | much more difficult | not at all | not at all fluent |

Table 3: Simplicity, adequacy and fluency scales used in the expert evaluation questionnaire (Q = question).

| Target texts | Simplicity | Adequacy | Fluency |
|---|---|---|---|
| RT | 4.25 | 3.75 | 3.75 |
| AS | 3.38 | 3.63 | 3.13 |
| MS | 4.38 | 4.25 | 4.25 |
| PE | 4.25 | 4.25 | 3.50 |

Table 4: Average simplicity, adequacy, and fluency ratings for each target text category (RT = reference target text, AS = automatically simplified text, MS = manually simplified text, PE = post-edited text).

consistently obtained the highest average ratings (i.e., 4.25, range of 3–5), while the automatically simplified texts were assigned the lowest average ratings (i.e., 3.63, range of 1–5, and 3.13, range of 2–4, respectively) (see Table 4).

The ratings of simplicity, adequacy, and fluency are consistent with the experts' comments in the evaluation questionnaire, in which seven out of eight experts reported finding the automatically simplified texts mostly not adequate. Simplification techniques were also considered to be only partially effective. On the other hand, the manually simplified texts were often reported as being "very good", even though it was also emphasized that they did not consistently comply with German Easy Language guidelines. As for the post-edited texts, most evaluators remarked on several simplicity as well as adequacy issues. Refer to Appendix A for examples of expert comments.

### 5.4. Discussion

The outcomes indicate that the automatic text simplification model under examination is not ready for deployment with or without post-editing, mainly due to weak simplification capabilities. As previously highlighted, the automatically simplified texts obtained, in fact, the lowest ratings across all evaluation metrics. In retrospect, it would have been beneficial to collect additional background data to identify the specific sets of Easy Language guidelines that experts commonly employed and referred to in their evaluations. Such data could have provided support for both the quantitative (i.e., ratings) and qualitative (i.e., comments) findings.

## 6. Overall Discussion

### 6.1. Advantage of Multi-stakeholder Involvement

Involving multiple stakeholders in ATS processes and assessments holds significant value for ensuring the ultimate quality and functionality of simplified texts. In our study, experts played a pivotal role by evaluating the adherence of texts to established guidelines, thereby offering critical insights into the simplicity, adequacy, and fluency of the simplified content (cf. Section 5). Conversely, post-editors contributed valuable feedback regarding productivity gains (cf. Section 4). Additionally, the perspectives of end users were indispensable for gauging the comprehensibility and acceptability of simplified texts in real-world contexts (cf. Section 3). It is crucial to acknowledge that linguistic complexity pertains to individual cognitive costs (Hansen-Schirra et al., 2020; Pallotti, 2015), and that text simplification efforts cater to highly diverse target groups. Hence, the active involvement of target audiences and the consideration of individual variability are paramount in optimizing the effectiveness and inclusivity of text simplification.

### 6.2. Advantage of Mixed-method Approaches

Mixed-method approaches offer several advantages in ATS studies. Given the inherent challenge of directly measuring reading comprehension behaviors, triangulating multiple proxies becomes imperative. The use of rating scales poses similar challenges, as interpretations may vary among participants (Stodden, 2021). Equally, qualitative findings regarding end users' perceptions of complex-

ity may diverge from quantitative metrics (Säuberli et al., 2024; Carrer, 2021; Benson-Goldberg et al., 2024), highlighting the importance of discerning discrepancies between perception and actual comprehension. Therefore, the adoption of mixed-method approaches not only enhances the robustness of research findings but also enables a more nuanced exploration of complex behaviors.

## 7. Conclusions

We conducted an extensive evaluation study involving end users, post-editors, and experts as stakeholders, and using a combination of quantitative/qualitative and objective/subjective methods. The results showed that there are differences in comprehensibility and perception of simplified texts between different user groups. We also found that qualitative surveys and behavioral observations can be essential in interpreting the results. These differences need to be accounted for in human evaluations of ATS models. Specifically, the following recommendations emerged from our experiments:

- Whenever possible, include target readers to assess comprehensibility.

- Do not rely solely on perceived quality ratings and assess the quality and functionality of the output as directly as possible, e.g., by measuring comprehension or post-editing effort.

- Collect qualitative data (e.g., through surveys or interviews) and behavioral measurements (e.g., while reading or post-editing) to support the interpretation of quantitative results.

- When collecting expert ratings, clearly define the Easy Language guidelines to be taken as a reference and ask the evaluators to specify which rules were not observed.

As ATS research begins to harness the new potential of large language models (Kew et al., 2023), future research should adopt a more human-centric and holistic approach to evaluation. We believe that this is essential for ensuring that the technological advancements yield tangible benefits for the end users of those technologies.

## 8. Acknowledgements

## 9. References

Suha S. Al-Thanyyan and Aqil M. Azmi. 2021. Automated text simplification: A survey. *ACM Comput. Surv.*, 54(2).

Oliver Alonzo, Jessica Trussell, Becca Dingman, and Matt Huenerfauth. 2021. Comparison of methods for evaluating complexity of simplified texts among deaf and hard-of-hearing adults at different literacy levels. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, CHI '21, Yokohama, Japan. Association for Computing Machinery.

Fernando Alva-Manchego, Carolina Scarton, and Lucia Specia. 2021. The (un)suitability of automatic evaluation metrics for text simplification. *Computational Linguistics*, 47(4):861–889.

Sergi Alvarez-Vidal and Antoni Oliver. 2023. Assessing mt with measures of pe effort. *Ampersand*, 11:100125.

Fabio Alves. 2003. *Triangulating Translation: Perspectives in process oriented research*. John Benjamins Publishing Company, Amsterdam.

Sofia Benson-Goldberg, Lori Geist, and Karen Erickson. 2024. Simplified COVID-19 guidance for adults with intellectual and developmental disabilities. *Journal of Applied Research in Intellectual Disabilities*, 37(3).

Bettina M. Bock and Sandra Pappert. 2023. *Leichte Sprache, Einfache Sprache, verständliche Sprache*. Narr Francke Attempto, Tübingen.

Ursula Bredel and Christiane Maaß. 2016. *Leichte Sprache. Theoretische Grundlagen. Orientierung für die Praxis*. Duden, Berlin.

Gilles Caporossi, Christophe Leblay, and Hakim Usoof. 2023. GenoGraphiX-Log version 2.0. User guide. *Les Cahiers Du GERAD G-2020-68*, pages 1–63.

Luisa Carrer. 2021. Translating into Easy Italian : an analysis of health-related texts and their impact on comprehension by people with intellectual disabilities. Thesis: Master, ZHAW Zürcher Hochschule für Angewandte Wissenschaften, Winterthur.

Edyta Charzyńska and Łukasz Jerzy Dębowski. 2015. Empirical verification of the polish formula of text difficulty. *Cognitive Studies*, 15:125–132.

Scott A Crossley, Hae Sung Yang, and Danielle S McNamara. 2014. What's so simple about simplified texts? a computational and psycholinguistic investigation of text comprehension and text processing. *Reading in a Foreign Language*, 26(1):92–113.

Inmaculada Fajardo, Vicenta Ávila, Antonio Ferrer, Gema Tavares, Marcos Gómez, and Ana Hernández. 2014. Easy-to-read texts for students with intellectual disability: linguistic factors affecting comprehension. *Journal of applied research in intellectual disabilities*, 27(3):212–225.

Christian Federmann. 2018. Appraise evaluation framework for machine translation. In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, pages 86–88, Santa Fe, New Mexico. Association for Computational Linguistics.

Jean-Paul Fox. 2010. *Bayesian Item Response Modeling: Theory and Applications*. Springer New York, New York, NY.

Sian Gooding, Yevgeni Berzak, Tony Mak, and Matt Sharifi. 2021. Predicting text readability from scrolling interactions. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 380–390, Online. Association for Computational Linguistics.

Natalia Grabar and Horacio Saggion. 2022. Evaluation of automatic text simplification: Where are we now, where should we go from here. In *Actes de la 29e Conférence sur le Traitement Automatique des Langues Naturelles. Volume 1 : conférence principale*, pages 453–463, Avignon, France. ATALA.

Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. Continuous measurement scales in human evaluation of machine translation. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 33–41, Sofia, Bulgaria. Association for Computational Linguistics.

Silvia Hansen-Schirra, Walter Bisang, Arne Nagels, Silke Gutermuth, Julia Fuchs, Liv Borghardt, Silvana Deilen, Anne-Kathrin Gros, Laura Schiffl, and Johanna Sommer. 2020. Intralingual translation into easy language – or how to reduce cognitive processing costs. In Silvia Hansen-Schirra and Christiane Maaß, editors, *Easy Language Research: Text and User Perspectives*, pages 197–225. Frank & Timme, Berlin.

Matt Huenerfauth, Lijun Feng, and Noémie Elhadad. 2009. Comparing evaluation techniques for text readability software for adults with intellectual disabilities. In *Proceedings of the 11th International ACM SIGACCESS Conference on Computers and Accessibility*, Assets '09, page 3–10, Pittsburgh, Pennsylvania, USA. Association for Computing Machinery.

Martin Kappus and Maureen Ehrensberger-Dow. 2020. The ergonomics of translation tools: understanding when less is actually more. *The Interpreter and Translator Trainer*, 14(4):386–404.

Tannon Kew, Alison Chi, Laura Vásquez-Rodríguez, Sweta Agrawal, Dennis Aumiller, Fernando Alva-Manchego, and Matthew Shardlow. 2023. BLESS: Benchmarking large language models on sentence simplification. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13291–13309, Singapore. Association for Computational Linguistics.

Hans P. Krings. 2001. *Repairing Texts: Empirical Investigations of Machine Translation Post-editing Processes*. Translation studies. Kent State University Press, Kent, Ohio and London.

Isabel Lacruz, Michael Denkowski, and Alon Lavie. 2014. Cognitive demand and cognitive effort in post-editing. In *Proceedings of the 11th Conference of the Association for Machine Translation in the Americas*, pages 73–84, Vancouver, Canada. Association for Machine Translation in the Americas.

Isabel Lacruz, Gregory M. Shreve, and Erik Angelone. 2012. Average pause ratio as an indicator of cognitive effort in post-editing: A case study. In *Workshop on Post-Editing Technology and Practice*, San Diego, California, USA. Association for Machine Translation in the Americas.

Gondy Leroy, James E Endicott, David Kauchak, Obay Mouradi, and Melissa Just. 2013. User evaluation of the effects of a text simplification algorithm using term familiarity on perception, understanding, learning, and information retention. *Journal of medical Internet research*, 15(7):e144.

Gondy Leroy, David Kauchak, Diane Haeger, and Douglas Spegman. 2022. Evaluation of an online text simplification editor using manual and automated metrics for perceived and actual text difficulty. *JAMIA Open*, 5(2):ooac044.

John Michael Linacre. 1989. *Many-faceted Rasch measurement*. Ph.D. thesis, The University of Chicago.

Christiane Maaß. 2020. *Easy Language–Plain Language–Easy Language Plus: Balancing comprehensibility and acceptability*. Frank & Timme, Berlin.

Jonathan Mallinson, Rico Sennrich, and Mirella Lapata. 2020. Zero-shot crosslingual sentence simplification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5109–5126, Online. Association for Computational Linguistics.

Louis Martin, Angela Fan, Éric de la Clergerie, Antoine Bordes, and Benoît Sagot. 2022. MUSS: Multilingual unsupervised sentence simplification by mining paraphrases. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1651–1664, Marseille, France. European Language Resources Association.

Gabriele Pallotti. 2015. A simple view of linguistic complexity. *Second Language Research*, 31(1):117–134.

Elissa Redmiles, Lisa Maszkiewicz, Emily Hwang, Dhruv Kuchhal, Everest Liu, Miraida Morales, Denis Peskov, Sudha Rao, Rock Stevens, Kristina Gligorić, Sean Kross, Michelle Mazurek, and Hal Daumé III. 2019. Comparing and developing tools to measure the readability of domain-specific texts. 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 4831–4842, Hong Kong, China. Association for Computational Linguistics.

Luz Rello, Ricardo Baeza-Yates, Stefan Bott, and Horacio Saggion. 2013a. Simplify or help? text simplification strategies for people with dyslexia. In *Proceedings of the 10th International Cross-Disciplinary Conference on Web Accessibility*, W4A '13, Rio de Janeiro, Brazil. Association for Computing Machinery.

Luz Rello, Ricardo Baeza-Yates, Laura Dempere-Marco, and Horacio Saggion. 2013b. Frequent words improve readability and short words improve understandability for people with dyslexia. In *INTERACT '13*, pages 203–219, Cape Town, South Africa. Springer.

Luz Rello, Susana Bautista, Ricardo Baeza-Yates, Pablo Gervás, Raquel Hervás, and Horacio Saggion. 2013c. One half or 50%? an eye-tracking study of number representation readability. In *INTERACT '13*, pages 229–245, Cape Town, South Africa. Springer.

Annette Rios, Nicolas Spring, Tannon Kew, Marek Kostrzewa, Andreas Säuberli, Mathias Müller, and Sarah Ebling. 2021. A new dataset and efficient baselines for document-level text simplification in German. In *Proceedings of the Third Workshop on New Frontiers in Summarization*,

pages 152–161, Online and in Dominican Republic. Association for Computational Linguistics.

Michael Ryan, Tarek Naous, and Wei Xu. 2023. Revisiting non-English text simplification: A unified multilingual benchmark. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4898–4927, Toronto, Canada. Association for Computational Linguistics.

Horacio Saggion, Sanja Štajner, Stefan Bott, Simon Mille, Luz Rello, and Biljana Drndarevic. 2015. Making it simplext: Implementation and evaluation of a text simplification system for spanish. *ACM Trans. Access. Comput.*, 6(4).

Sanja Štajner and Sergiu Nisioi. 2018. A detailed evaluation of neural sequence-to-sequence models for in-domain and cross-domain text simplification. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Regina Stodden. 2021. When the scale is unclear - analysis of the interpretation of rating scales in human evaluation of text simplification. In *Proceedings of the First Workshop on Current Trends in Text Simplification (CTTS 2021) co-located with the 37th Conference of the Spanish Society for Natural Language Processing (SEPLN2021), Online (initially located in Málaga, Spain), September 21st, 2021*, volume 2944 of *CEUR Workshop Proceedings*. CEUR-WS.org.

Andreas Säuberli, Franz Holzknecht, Patrick Haller, Silvana Deilen, Laura Schiffl, Silvia Hansen-Schirra, and Sarah Ebling. 2024. Digital comprehensibility assessment of simplified texts among persons with intellectual disabilities.

Sanja Štajner. 2021. Automatic text simplification for social good: Progress and challenges. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2637–2652, Online. Association for Computational Linguistics.

# A.  Expert comments

The following is a selection of comments from the expert evaluation. Comments were translated into English from German.

## A.1.  Automatically simplified text (AS)

7 out 8 experts reported finding the AS texts mostly not adequate.

- "The TT contains some 'information vagueness'"
- "In terms of content, these are two different texts"
- "The simplification is mainly achieved through a different text structure"
- "This text is not in Easy Language. At most, it is a shortened version of the original text"

## A.2.  Manually simplified text (MS)

Often reported as being "very good".

- "very good"
- "One of the better texts in this questionnaire"
- "Again, this text is not in Easy Language, although it does comply with many of the rules"
- "The rules are not consistently adhered to"

## A.3.  Reference text (RT)

More in line with Plain Language properties.

- "One of the better texts here and much easier to understand for the target group"
- "The target text combines two different language levels"
- "Here too: Target text is no Easy Language" / "This text is not in Easy Language"

## A.4.  Post-edited text (PE)

Experts reported simplicity/adequacy issues.

- "The target text is not written in Easy Language: Several rules are not observed"
- "Significant reduction in content"
- "The TT still contains some difficult words"

# Decoding the Metrics Maze: Navigating the Landscape of Conversational Question Answering System Evaluation in Procedural Tasks

## Alexander Frummet and David Elsweiler

University of Regensburg

{alexander.frummet, david.elsweiler}@ur.de

## Abstract

Conversational systems are widely used for various tasks, from answering general questions to domain-specific procedural tasks, such as cooking. While the effectiveness of metrics for evaluating general question answering (QA) tasks has been extensively studied, the evaluation of procedural QA remains a challenge as we do not know what answer types users prefer in such tasks. Existing studies on metrics evaluation often focus on general QA tasks and typically limit assessments to one answer type, such as short, SQuAD-like responses or longer passages. This research aims to achieve two objectives. Firstly, it seeks to identify the desired traits of conversational QA systems in procedural tasks, particularly in the context of cooking (RQ1). Second, it assesses how commonly used conversational QA metrics align with these traits and perform across various categories of correct and incorrect answers (RQ2). Our findings reveal that users generally favour concise conversational responses, except in time-sensitive scenarios where brief, clear answers hold more value (e.g. when heating in oil). While metrics effectively identify inaccuracies in short responses, several commonly employed metrics tend to assign higher scores to incorrect conversational answers when compared to correct ones. We provide a selection of metrics that reliably detect correct and incorrect information in short and conversational answers.

**Keywords:** metrics, conversational search, question answering

## 1. Introduction

Conversational systems are frequently used for a variety of tasks, such as setting timers, getting the weather forecast for the day, or retrieving factual information from the web. For such general question answering (QA) tasks, users can ask questions such as "What is the capital of Germany?". Responses from conversational agents can vary, from short, concise answers, e.g., "Berlin", to more conversational responses, such as "The capital of Germany is Berlin." However, the accuracy of these different answer types cannot be guaranteed. Therefore, there is a crucial need for reliable metrics to evaluate the effectiveness of such conversational systems. For general QA tasks, many studies have explored the effectiveness of commonly used metrics, spanning *word overlap-based metrics*, e.g., BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), F1 Score, and Exact Match, measuring the degree of overlap between words in the ground truth answer and those generated by a QA model, to *embedding-based metrics*, such as Semantic Answer Similarity (SAS) (Risch et al., 2021) and BERTScore (Zhang et al., 2020), which measure semantically equivalent responses.

To determine the utility of these metrics, studies commonly examine their alignment with user preferences. For instance, numerous studies have examined the efficacy of word overlap-based metrics, and have concluded that METEOR exhibits

the strongest correlation with human evaluations of QA model outputs, in comparison to metrics, such as ROUGE and BLEU (Blagec et al., 2020; Nema and Khapra, 2018; Chen et al., 2019). Others have investigated the performance of embedding-based metrics, including SAS and BERTScore. Chen et al. (2019) found that while BERTScore is superior at capturing semantic information, it does not correlate as strongly with human assessments of how closely a model prediction matches the ground truth. Additionally, SAS and BERTScore struggle with spatial awareness, numbers, and conversions (Mustafazade et al., 2022). In one particularly thorough analysis by Liu et al. (2021), conversational search evaluation metrics were evaluated from three different perspectives: reliability, fidelity, and intuitiveness. Based on their analysis, METEOR and BERTScore were determined to be the most reliable, whereas METEOR and BLEU were found to be the most intuitive.

These studies primarily assess metrics for general, domain-agnostic QA using open-domain datasets such as MSDialog or Wizard of Wikipedia (Liu et al., 2021). However, for domain-specific procedural QA tasks, such as cooking and DIY, where conversational agents are increasingly popular[1], such meta-evaluations fall short. While these domains share similarities with general QA, we lack

---

[1]see cooking/DIY focus in Alexa Task-Bot challenge: https://www.amazon.science/alexa-prize/taskbot-challenge

a deep understanding of the unique challenges in procedural QA conversations, making evaluation difficult. For example, consider the recipe step "Add apples, oranges, and sugar to a large pitcher and muddle with a muddler or wooden spoon for 45 seconds". When asked, "How much sugar do I need to add?" the answer can be as concise as "50g" or more conversational, e.g., "You need to add 50g of sugar". These responses differ in the level of context provided about the cooking process. "50g" lacks context, while the other response specifies that "50g" refers to sugar. In-situ studies with voice assistants have shown that users have a preference for agents that provide clarifications (Luger and Sellen, 2016) and in the cooking domain, an analysis of human-human dialogues revealed that people often seek reassurances regarding the answers they receive from an assistant (Frummet et al., 2022).

While these studies hint that answers clarifying context may be preferred, we do not yet know what constitutes a good answer in procedural tasks, such as cooking. Since we lack knowledge of user answer preferences, selecting a reliable metric to evaluate answer correctness in procedural QA tasks remains challenging.

This differs from "general" non-procedural QA tasks where studies evaluate metric effectiveness using responses of the same type, whether short, SQuAD-like answers (Mustafazade et al., 2022; Nema and Khapra, 2018; Bulian et al., 2022), conversational, sentence-length answers (Shi et al., 2023; Siblini et al., 2021) or long, paragraph-like answers (Xu et al., 2023). While these studies provide insights into reliable metrics for non-procedural scenarios, it is unclear if these findings extend to procedural tasks. It is uncertain if the metrics used there remain reliable, valid, and applicable when various answer formulations are possible, as is the case in procedural tasks.

To assess conversational agents effectively for procedural tasks, we need to understand 1) which answers are preferred by users and 2) whether the metrics traditionally employed in conversational QA yield reliable results for evaluating conversational systems in procedural tasks.

Having reliable evaluation results is crucial for the success of most procedural tasks, as having the correct information, such as quantities and next steps, is vital to successfully completing the associated task (Frummet et al., 2022). Just as with human-generated responses, multiple answer formulations are possible, and metrics must account for these variations. Additionally, they must be sensitive to incorrect aspects of answers, given that large language models sometimes "just make stuff up" (Shah and Bender, 2022). Users tend to trust these models because they mimic human language

(Araujo, 2018; Dinan et al., 2021).

This study aims to achieve two objectives within the context of procedural assistance tasks. Firstly, it seeks to identify the desirable traits of QA systems for humans (RQ1). Secondly, it aims to analyse how commonly used metrics in conversational QA reflect these traits and vary for different categories of correct and incorrect answers (RQ2).

## 2. Methodology

This section outlines the methods and resources used to address the research questions, including the dataset used, the various answer types and metrics evaluated, and a user study to complement our system-based analyses.

### 2.1. Dataset

In this paper, we target cooking-related procedural assistance tasks. Existing datasets, such as Cook-Dial (Jiang et al., 2022) and Wizard of Tasks (Choi et al., 2022), require post-processing to meet our needs. For example, Wizard of Tasks lacks grounding for conversational answers to specific parts of the recipe, making it challenging to evaluate different answer types within the recipe context. To address these limitations, we have created a new dataset tailored to our study's requirements. As a basis for our experiments, we use 298 randomly selected questions and answers [2] from a conversational cooking QA dataset (QookA) (Frummet and Elsweiler, 2024). The full dataset contains 1268 pairs of question-answer, where the questions are expressed in written natural language after being transcribed from spoken questions gathered from 95 participants ($M_{age} = 35$, $73\%$ female, $26\%$ male, $1\%$ diverse) who followed a recipe from SeriousEats[3] in a simulated cooking scenario. Each question maps to an information need type from the taxonomy presented in Frummet et al. (2022) and has an associated answer derived from the appropriate recipe. Examples can be seen in the first two columns of Table 1.

### 2.2. Answer Formulations

Our study examines how metrics relate to various styles of answer. To this end, we devised three distinct answer types, each of which can be considered as "correct". First, *Classic QA* answers adhere to the extractive, span-based SQuAD format (Rajpurkar et al., 2018) and contain only the relevant information that is needed to answer a question.

| QookA Dataset | Classic QA | | Conv. QA | | Passage |
|---|---|---|---|---|---|
| Question | Answer | HA1 | HA2 | GPT-3 | |
| What do I add after the mustard and cheese? | eggs. | eggs. | After mustard and cheese are combined, you need to add eggs one at a time. | After adding the mustard and cheese, you should add the eggs, one at a time. | After adding the mustard and cheese, beat the mixture with a wooden spoon until homogeneous. Then add eggs one at a time, ... | To Finish by Hand: Remove pot from heat. Add mustard and cheese and beat with wooden spoon until ... |
| How much fat should I reserve? | 1 cup. | 1 cup. | After skimming off the fat, 1 cup should be left behind. | According to the recipe you should reserve about 1 cup of the solid fat | Reserve 1 cup of fat when skimming the sauce. | Carefully skim off most of the fat, leaving behind about 1 cup total. (For more precise measurement … |
| How long will it take for the coconut fat to separate? | about 4 minutes. | about 4 minutes. | It will take about 4 minutes for the coconut fat to separate. | Until coconut fat separates and solids start to sizzle, about 4 minutes. | It will take approximately 4 minutes for the coconut fat to separate when heated over medium-high heat and stirred frequently. | Open 1 can of coconut milk and carefully spoon off about 3 tablespoons of the solid fat from the surface and place it in a large Dutch oven... |

Table 1: Example entries from the dataset used. The left-most two columns are from the QookA dataset. The remaining columns represent the answer types studied in this work. HA1 and HA2 denote the two human annotators.

Second, we investigate answers that follow the style of *conversational QA* datasets, such as QuAC (Choi et al., 2018), CoQA (Reddy et al., 2019) and Wizard of Tasks (Choi et al., 2022). To establish these answers, two human annotators were asked to formulate what they would consider to be "ideal" responses given a question and the classic QA answer, as well as the step text from the QookA dataset. They were free to formulate the answer in any way they wished, and did so independently. As an automated comparison, we provided the same instructions and information to GPT-3 DaVinci-3 as a prompt:

Question: ⟨*Question*⟩
Answer: ⟨*Answer*⟩
Context: ⟨*Recipe Step Text*⟩
Rephrased Answer:

The authors checked these manually to establish that they were still "correct".

Last, we investigate passages that contain the answer and the surrounding context. These were attained by identifying the recipe step that included the pertinent information. A passage-based answer is appropriate since answers of this type are evaluated in conversational information retrieval (IR) assessment frameworks, such as CAsT (Dalton et al., 2020), and is a plausible information unit to present to users since past research has revealed that cooking assistant users value contextually embedded answers (Frummet et al., 2019, 2022). Examples answers of all three types can be found in Table 1.

## 2.3. Incorrect answers

To obtain incorrect versions of the same classes of answer we leveraged GPT-3 to generate responses that were factually inaccurate using the following prompt:

Question: ⟨*Question*⟩
Context: ⟨*Recipe Step Text*⟩
Correct Answer: ⟨*Correct Answer (either ClassicQA/Conv. QA)*⟩
Factually wrong answer:

GPT-3 was provided with the question, the correct answer (either Classic QA or Conversational QA form) and the corresponding recipe step. In the case of Conversational QA, we randomly chose one of the two human annotators and used their answer as the ground truth. As an example, for the second row in Table 1, the incorrect Classic QA answer was "a tablespoon" and the incorrect Conversational QA answer was "You should reserve 2 cups of fat".

To derive incorrect passage formulations, we opted to randomly select another passage from our dataset to replicate the potential for an improperly retrieved passage that may arise during a TREC CAsT experiment.

## 2.4. Metrics studied

We study the most commonly applied QA metrics of the classes outlined in the introduction: Word-overlap based and embbeding-based metrics. Specifically, we evaluate ROUGE (Lin, 2004), BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), Exact Match and F1 for the word-overlap domain. From the embedding domain, we study Semantic Answer Similarity (Risch et al., 2021) and BERTScore F1 (Zhang et al., 2020).

To calculate the metric, we selected one of our answers as the prediction and compared it against all other answers. For example, if testing a classic QA answer (i.e., prediction) then the equivalent HA1, HA2, GPT-3, and Passage answers served as the references/ground truth to compute the metric values. We utilised Huggingface's evaluation library [4] to calculate the ROUGE, BERTScore, METEOR, BLEU, and Exact Match scores, while FARM's evaluation library [5] was used to compute F1 scores. Additionally, we customised the script from the SCAI-QReCC-22 shared task to determine the Semantic Answer Similarity (SAS).

## 2.5. User perception of correct answers

We conducted an online experiment with a within-groups design to determine how participants perceive answers of different types. In this experiment, each participant evaluated the appropriateness of five separate answers (one for each type, presented in a random order) for questions chosen

from our dataset. As illustrated in Figure 1, ratings were provided on a Likert scale ranging from 1 to 5. To provide the necessary context, each answer was accompanied by the corresponding question and recipe step. To learn why our participants preferred certain answer types over others, we requested that they provide an explanation for their ratings.

The experiment was designed to resemble an interaction with a conversational assistant in a kitchen. To this end, we employed Google's Text-to-Speech API to convert the answers to audio files, which participants then listened to. Attention checks were used to confirm that participants had indeed listened to the answers. As an attention check, one of the answer audio files illustrated in Figure 1 contained the following instruction: "Please click the left circle and write the answer to six multiplied by four into the text field below."

Study participants were recruited via Prolific. According to the power analysis performed using G*Power (Faul et al., 2007), a total of 32 individuals were needed to achieve the required statistical power for conducting an ANOVA test with repeated measures [6]. All participants were recruited from the UK as we selected a British accent for the audio answer files. 53.13% of our participants were female, 46.88% male, with most being between 25 and 34 years old (37.5%). Since our experiment is in the cooking domain, we wanted to know how much people enjoy cooking on a scale from 1 to 5. The people in our study generally enjoy cooking ($M = 3.66, SD = 1.32$).

## 3. Results

This section presents our findings. Sections 3.1 and 3.2 provide an insight into the human perspective (RQ1) by examining the answers provided by human annotators and reporting the outcomes of our online study. Meanwhile, Sections 3.4 and 3.5 shed light on the metrics' ability to differentiate between correct and incorrect answers, and how they reflect the characteristics users desire (RQ2).

## 3.1. Evaluating Human Provided Answers

In a first step, we analyse the "ideal" answers provided by two human annotators to understand the conveyed information and the methods of communication. An initial observation suggests that these answers are considerably lengthier than Classic QA answers (which have a mean word count of $\bar{x} = 3.69$), but shorter than passage answers ($\bar{x} = 71.67$). On average, Human annotator 1 (HA1)

---

[4] https://huggingface.co/evaluate-metric
[5] https://farm.deepset.ai/_modules/farm/eval.html

---

[6] desired Power: 0.8; effect size $f = 0.25$; significance threshold $\alpha = 0.05$; num. of groups: 1; num. of measurements: 5; Corr. among rep. measures: 0.2

**Context if needed**

**Recipe Title:** Crispy Kung Pao Tofu Recipe

**Current Step:** Heat oil in a wok to 350°F. Whisk together 1/2 cup cornstarch, flour, baking powder, and 1 teaspoon kosher salt. Add water and vodka and whisk until a smooth batter is formed, adding up to 2 tablespoons additional water if batter is too thick. It should have the consistency of thin paint and fall off of the whisk in thin ribbons that instantly disappear as they hit the surface of the batter in the bowl.

Below you will find a question and the corresponding answer provided by a digital cooking assistant.

**Question:** what should I use to mix the cornstarch flour baking powder and salt?

**Answer:**

▶ 0:00 / 0:00 🔊 ⋮

**How appropriate is the provided answer for the question shown?** (The yellow box above may be helpful.)

inappropriate ○ ○ ○ ○ ○ appropriate

What made the snippet appropriate or inappropriate? *

PREVIOUS                                                                 NEXT

Figure 1: Screenshot of our online experiment tool.

provided responses containing $\bar{x} = 12.40$ words, while Human annotator 2 (HA2) had answers averaging $\bar{x} = 9.78$ words.

Examining the answers more closely reveals that both human annotators chose their formulations carefully. Many of the answers were phrased, such that the user would be reassured that the question was correctly understood. The first entry in Table 1 exemplifies that both HA1 and HA2 embed their answer in the conversational and step context. Both answers start with "After adding the mustard and cheese...". This linguistic pattern is used to reassure the user that their question has been understood correctly, a pattern observed in naturalistic cooking QA investigations (Frummet et al., 2022). Rather than simply providing times, e.g., "about 4 minutes", as the classic QA answer does, both HAs provide contextual information to provide cues to aid understanding e.g., "It will take

about 4 minutes for the coconut fat to separate". Such techniques were common as is reflected in the jaccard similarity between HAs and questions (20% HA1, 11.58% HA2) and HAs and passage (13.15% HA1, 12.12% HA2).

## 3.2. Establishing User Answer Preferences

Participants rated conversational QA answers higher (GPT-3: $M = 4.53, SD = .67$, HA1: $M = 4.34, SD = 1.07$, HA2: $M = 4.22, SD = 1.10$) than classic QA answers ($M = 3.97, SD = 1.36$) and passage answers ($M = 3.00, SD = 1.34$). An ANOVA with repeated measures showed a significant difference ($F(4, 124) = 9.14, p < .001$). Bonferroni-adjusted post-hoc t-tests revealed that passage answers were rated significantly lower than GPT-3 answers ($p < .001$), HA1 answers

| Answer Type | ROUGE1 | ROUGE2 | ROUGEL | ROUGELSum | BERTScore | SAS | BLEU | METEOR | EM | F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| Cor. Classic QA | 44.03 | 30.64 | 38.67 | 38.67 | 90.36 | 68.53 | 21.89 | 58.63 | 0.34 | 43.29 |
| Inc. Classic QA | 13.50 | 1.88 | 13.35 | 13.36 | 92.59 | 42.83 | 1.80 | 26.25 | 0.00 | 12.62 |
| Cor. Conv. QA | 58.37 | 39.95 | 54.68 | 54.68 | 93.41 | 78.85 | 34.63 | 61.99 | 1.51 | 57.48 |
| Inc. Conv. QA | 59.29 | 43.87 | 58.39 | 58.39 | 94.42 | 56.67 | 38.86 | 56.52 | 0.34 | 56.66 |
| Cor. Passage | 45.24 | 25.49 | 45.02 | 45.02 | 90.54 | 67.45 | 15.32 | 32.87 | 2.01 | 49.16 |
| Inc. Passage | 22.01 | 4.66 | 15.56 | 15.56 | 84.67 | 44.43 | 2.63 | 21.07 | 1.34 | 19.76 |

Table 2: Metric Experiment results (in percent) grouped by Answer Type. The results for Conv. QA answer types are the mean scores for HA1, HA2 and GPT-3, which show little variation.

($p < .001$) and HA2 answers ($p < .01$).

The results indicate a preference for conversational QA answers compared to other types, which is supported by the justifications provided with the ratings.

## 3.3. Understanding User Answer Preferences

To gain insights into participants' preferences for specific answers, we conducted a qualitative analysis of the explanations they provided for their ratings.

### 3.3.1. Passage

In line with the ratings discussed in Section 3.2, the majority of feedback related to passage-style answers was negative (84.38%). Participants often criticised these answers for being overly lengthy, containing "too much information", and lacking a "direct" and "specific" response. Consequently, they found these responses "confusing" and "unhelpful". However, some participants did find them "totally appropriate" without providing further elaboration.

### 3.3.2. Classic QA

Conversely, explanations for classic QA responses were more balanced, with 56.25% of the feedback being positive and 43.75% negative. Participants appreciated the brevity and clarity of these responses, finding them "short", "to the point", "clear", and "concise". Some mentioned that they particularly liked these answers in situations requiring "quick" reactions, for example, when they are "in the midst of heating oil". However, in contexts not demanding immediate responses, participants critiqued the lack of detail (e.g., "relevant but too brief", "could be more exact and informative", "I think the AI should be clearer") and called for more contextual information to fully comprehend the answer within the cooking process. For instance, a participant felt that from the given answer it was "unclear when to add eggs". Another individual said: "It correctly said stand mixer but could have added that it needs a paddle attachment". Others suggested providing

"extra information", such as instructions on how "to score with a knife".

### 3.3.3. Conversational QA

Conversational answers generated by HA1, HA2, and GPT-3 received the most positive feedback, with an average of 73.94% of answers rated positively and 21.88% rated negatively. Participants favored these responses for being straightforward, clear, and appropriately detailed. Unlike classic QA, conversational answers provided "a good amount of detail" and offered extra context that "helps to clarify what part of the recipe it is referring to at the same time as getting the answer". Participants noted that these answers helped them plan ahead in the cooking process, (e.g., "The answer told me to put the ingredients in a bowl, however it also went beyond and spoke about the type of mixer.", "useful to know when to sprinkle the parmesan", "vocalised in chronological order the steps necessary to complete this section of the recipe").

The few negative comments suggested some answers could be shorter for conciseness (e.g., "'In a large skillet' would've been enough and more concise").

### 3.3.4. Summary

In conclusion, our analysis of user preferences for different answer types reveals the following key insights:

- **Passage**: Users predominantly criticise lengthy, unclear passage-style answers, often finding them confusing and unhelpful.

- **Classic QA**: Classic QA answers are appreciated for their brevity, particularly in situations requiring quick responses, but some users call for more detail to fully understand the context.

- **Conversational QA**: Conversational answers generated the most positive feedback due to their clarity, detail, and suitability for planning the cooking process. However, a few users suggested that some responses could

be made more concise for improved user experience.

### 3.4. Understanding variance in metric scores across type of correct answer

The effectiveness of the metrics were assessed by grouping the results by type of answer (Classic QA, Conversational QA, Passage). As indicated in Table 2, the exact match metric yielded extremely poor results across all conditions (i.e., $< 2.5\%$). Both the conventional machine translation metrics, BLEU, ROUGE, and METEOR, and the commonly used F-Measure also achieved low scores for correct answers. Of these, however, METEOR provided the highest scores, which aligns with past research indicating that METEOR is the most robust among the common Machine Translation metrics (see, for example, Chen et al. (2019); Blagec et al. (2020); Nema and Khapra (2018)). The embedding-based metrics BERTScore and SAS, in contrast, yield much higher scores than all other metrics for correct answers with BERTScore providing the highest of all. Contrasting the metric scores across the three classes of answer reveals that all metrics provide higher scores for the user-preferred conversational answers.

### 3.5. Understanding variance in metric scores across type of incorrect answer

Our results show that all metrics provide lower scores when we provide factually incorrect classic QA and passage style answers. This is an expected and desired outcome. However, as illustrated in Table 2, most metrics yield higher scores when an incorrect answer in the conversational QA format is provided. BERTScore even yields the highest scores overall ($> 94\%$). This indicates that BERTScore is not a suitable metric for dealing with incorrect information. The only exceptions are the SAS and METEOR metrics which decrease for incorrect answers. The decreasing METEOR metric score, again, evidences its robustness as pointed out in Chen et al. (2019) and Liu et al. (2021).

We performed a Kruskal-Wallis test to determine if the differences in ConvQA answer type are statistically significant. The different metrics results served as dependent variables. Our independent variable is correct/incorrect answer. A Posthoc Dunn's test with bonferroni-adjusted p-values revealed that ROUGE2/L/LSum ($p < 0.05$) and BERTScore ($p < 0.001$) achieved significant higher scores for incorrect answers. METEOR ($p < 0.01$) and SAS ($p < 0.001$) yielded significantly higher results for correct answers.

## 4. Discussion

In this work, we have tried to "decode the metrics maze" by evaluating popular question answering metrics in the light of two research questions.

### 4.1. User Preferences for Answer Types

In RQ1, we studied the desired traits of procedural QA systems for users and whether all types of correct answers were equally preferred. Our findings from the qualitative analysis provide some hints that user preferences for answer types may vary depending on the specific task they perform. In situations requiring quick responses, such as boiling or heating something in oil, users tend to favour short and to the point answers. During most stages of the cooking process, however, users prefer concise, contextual conversational answers helping them to plan the cooking process. It's evident that not all correct answers are viewed equally by users, aligning with observations from previous naturalistic research on conversational cooking QA scenarios (Frummet et al., 2019, 2022, 2024). Consequently, for procedural conversational question answering, such as cooking, metrics should reliably distinguish between correct and incorrect answers for both short (=Classic QA) and conversational responses.

### 4.2. Metrics Performance in Procedural Tasks

In RQ2, we examined how commonly used metrics in conversational QA align with user preferences from RQ1 and their performance across various answer categories in a cooking QA scenario. Users generally favour conversational answers but appreciate short, concise responses in specific situations. To meet this need, metrics should effectively distinguish between correct and incorrect answers in both response styles.

Our findings, presented in sections 3.4 and 3.5, revealed that, except for SAS, METEOR and F1, all metrics exhibited higher scores for incorrect answers compared to correct ones in the case of conversational responses. This trend was consistent for metrics commonly employed in machine learning tasks, including ROUGE, BLEU, and BERTScore, with METEOR being the only exception. Hanna and Bojar (2021) provided an explanation for the generally high BERTScore values stating that "BERTScore fails to assign low scores when a bad candidate sentence has high lexical overlap with the reference in terms of content words" (Hanna and Bojar, 2021, p. 515). This phenomenon is attributed to the lexical similarity between correct and incorrect answers in our study, where only a few words differ. Both Blagec et al.

(2020) and Chen et al. (2019) have argued that machine translation metrics such as ROUGE and BLEU are unsuitable for question answering tasks as they struggle to identify incorrect information due to their n-gram-based approach (Blagec et al., 2020). However, Blagec et al. (2020) noted that METEOR is good at capturing semantic differences.

Our results suggest that commonly used conversational QA metrics may not accurately evaluate the correctness of information. Instead, a comprehensive suite of metrics is needed to assess answer accuracy. For short answers, F1, METEOR and SAS are dependable choices for distinguishing between correct and incorrect information. In the case of conversational responses, METEOR and SAS are better choices than F1 in capturing these distinctions more effectively. The choice of metrics should align with the answer type preferred by the user in the current cooking task (see Section 4.1). For short answers, the appropriate suite includes F1, METEOR, and SAS, while for conversational answers, it consists of METEOR and SAS.

## 5. Conclusion and Future Work

Currently, conversational assistant systems are being evaluated with diverse and, as we have shown here, potentially inappropriate metrics. We suggest using a suite of metrics to accurately assess the effectiveness of such systems. The choice of metrics should be tailored to the specific task and how users respond within that task.

Although we discovered hints that users tend to favour brief answers during time-sensitive situations (e.g., heating in oil) and conversational answers in most other cooking stages, the generalisability of these preferences remains uncertain. Future work should investigate these critical scenarios further. Moreover, forthcoming research should focus on developing robust evaluation methods for handling inaccuracies. While our study on a cooking QA dataset underscores this challenge, further confirmation using other datasets such as Wizard of Tasks (Choi et al., 2022) and CookDial (Jiang et al., 2022) is needed.

## 6. Bibliographical References

Theo Araujo. 2018. Living up to the chatbot hype: The influence of anthropomorphic design cues and communicative agency framing on conversational agent and company perceptions. *Computers in Human Behavior*, 85:183–189.

Satanjeev Banerjee and Alon Lavie. 2005. ME-TEOR: 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. Association for Computational Linguistics.

Kathrin Blagec, Georg Dorffner, Milad Moradi, and Matthias Samwald. 2020. A critical analysis of metrics used for measuring progress in artificial intelligence. *CoRR*, abs/2008.02577.

Jannis Bulian, Christian Buck, Wojciech Gajewski, Benjamin Börschinger, and Tal Schuster. 2022. Tomayto, tomahto. beyond token-level answer equivalence for question answering evaluation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 291–305, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Anthony Chen, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. Evaluating question answering evaluation. In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 119–124, Hong Kong, China. Association for Computational Linguistics.

Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wen-tau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. Quac: Question answering in context. *arXiv preprint arXiv:1808.07036*.

Jason Ingyu Choi, Saar Kuzi, Nikhita Vedula, Jie Zhao, Giuseppe Castellucci, Marcus Collins, Shervin Malmasi, Oleg Rokhlenko, and Eugene Agichtein. 2022. Wizard of tasks: A novel conversational dataset for solving real-world tasks in conversational settings. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3514–3529, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Jeffrey Dalton, Chenyan Xiong, Vaibhav Kumar, and Jamie Callan. 2020. Cast-19: A dataset for conversational information seeking. *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*.

Emily Dinan, Gavin Abercrombie, A Stevie Bergman, Shannon Spruit, Dirk Hovy, Y-Lan Boureau, and Verena Rieser. 2021. Anticipating safety issues in e2e conversational ai: Framework and tooling. *arXiv preprint arXiv:2107.03451*.

Franz Faul, Edgar Erdfelder, Albert-Georg Lang, and Axel Buchner. 2007. G*power 3: A flexible statistical power analysis program for the social,

behavioral, and biomedical sciences. *Behavior Research Methods*, 39(2):175–191.

Alexander Frummet and David Elsweiler. 2024. Qooka: A cooking question answering dataset. In *Proceedings of the 2024 ACM SIGIR Conference on Human Information Interaction and Retrieval, CHIIR 2024, Sheffield, United Kingdom, March 10-14, 2024*, pages 406–410. ACM.

Alexander Frummet, David Elsweiler, and Bernd Ludwig. 2019. Detecting domain-specific information needs in conversational search dialogues. In *Proceedings of the 3rd Workshop on Natural Language for Artificial Intelligence co-located with the 18th International Conference of the Italian Association for Artificial Intelligence (AIIA 2019), Rende, Italy, November 19th-22nd, 2019*, volume 2521 of *CEUR Workshop Proceedings*. CEUR-WS.org.

Alexander Frummet, David Elsweiler, and Bernd Ludwig. 2022. "what can I cook with these ingredients?" - understanding cooking-related information needs in conversational search. *ACM Trans. Inf. Syst.*, 40(4):81:1–81:32.

Alexander Frummet, Alessandro Speggiorin, David Elsweiler, Anton Leuski, and Jeff Dalton. 2024. Cooking with conversation: Enhancing user engagement and learning with a knowledge-enhancing assistant. *ACM Trans. Inf. Syst.* Just Accepted.

Michael Hanna and Ondřej Bojar. 2021. A fine-grained analysis of BERTScore. In *Proceedings of the Sixth Conference on Machine Translation*, pages 507–517, Online. Association for Computational Linguistics.

Yiwei Jiang, Klim Zaporojets, Johannes Deleu, Thomas Demeester, and Chris Develder. 2022. Cookdial: A dataset for task-oriented dialogs grounded in procedural documents. *Applied Intelligence*, 53(4):4748–4766.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Zeyang Liu, Ke Zhou, and Max L. Wilson. 2021. Meta-evaluation of conversational search evaluation metrics. *ACM Trans. Inf. Syst.*, 39(4).

Ewa Luger and Abigail Sellen. 2016. "like having a really bad pa": The gulf between user expectation and experience of conversational agents. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, CHI '16, page 5286–5297, New York, NY, USA. Association for Computing Machinery.

Farida Mustafazade, Peter Ebbinghaus, and Seth Darren. 2022. Evaluation of semantic answer similarity metrics. *International Journal on Natural Language Computing*, 11:15.

Preksha Nema and Mitesh M. Khapra. 2018. Towards a better metric for evaluating question generation systems. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3950–3959, Brussels, Belgium. Association for Computational Linguistics.

Kishore Papineni, 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 '02, page 311–318, USA. Association for Computational Linguistics.

Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics.

Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. CoQA: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266.

Julian Risch, Timo Möller, Julian Gutsch, and Malte Pietsch. 2021. Semantic answer similarity for evaluating question answering models. In *Proceedings of the 3rd Workshop on Machine Reading for Question Answering*, pages 149–157, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Chirag Shah and Emily M Bender. 2022. Situating search. In *ACM SIGIR Conference on Human Information Interaction and Retrieval*, pages 221–232.

Zhengliang Shi, Weiwei Sun, Shuo Zhang, Zhen Zhang, Pengjie Ren, and Zhaochun Ren. 2023. RADE: Reference-assisted dialogue evaluation for open-domain dialogue. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12856–12875, Toronto, Canada. Association for Computational Linguistics.

Wissam Siblini, Baris Sayil, and Yacine Kessaci. 2021. Towards a more robust evaluation for conversational question answering. In *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 1028–1034, Online. Association for Computational Linguistics.

Fangyuan Xu, Yixiao Song, Mohit Iyyer, and Eunsol Choi. 2023. A critical evaluation of evaluations for long-form question answering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3225–3245, Toronto, Canada. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with BERT. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

# The 2024 ReproNLP Shared Task on Reproducibility of Evaluations in NLP: Overview and Results

**Anya Belz, Craig Thomson**

ADAPT Research Centre, Dublin City University

Dublin, Ireland

{anya.belz,craig.thomson}@dcu.ie

## Abstract

This paper presents an overview of, and the results from, the 2024 Shared Task on Reproducibility of Evaluations in NLP (ReproNLP'24), following on from three previous shared tasks on reproducibility of evaluations in NLP, ReproNLP'23, ReproGen'22 and ReproGen'21. This shared task series forms part of an ongoing research programme designed to develop theory and practice of reproducibility assessment in NLP and machine learning, against a backdrop of increasing recognition of the importance of reproducibility across the two fields. We describe the ReproNLP'24 shared task, summarise results from the reproduction studies submitted, and provide additional comparative analysis of their results.

**Keywords:** Reproducibility, Shared Task, Evaluation.

## 1. Introduction

Reproducibility continues to be a problem in search of a solution in the Natural Language Processing (NLP) field (Belz et al., 2021a, 2023). We still do not understand well enough what makes system evaluations, both human and metric-based, easier or harder to reproduce, while a growing number of reproduction studies have revealed alarmingly poor degrees of reproducibility and numerous issues with current evaluation practices (Belz et al., 2023).

The aim of this fifth reproduction-focused shared task in NLP, following REPROLANG'20 (Branco et al., 2020), ReproGen'21 (Belz et al., 2021b), ReproGen'22 (Belz et al., 2022), and ReproNLP'23 (Belz and Thomson, 2023), is generally to continue to add to the body of reproduction studies in NLP and machine learning (ML), and more specifically, to produce and analyse multiple reproductions of shared original evaluations, thereby creating more reliable reproducibility results for individual evaluations and evaluation methods, given that the evidence is that multiple reproductions rarely produce the same reproducibility results.

The 19 new reproduction studies that make up ReproNLP'24 add a good number of further data points available for investigating reproducibility, and help to continue identifying properties of evaluations that are associated with better reproducibility.

We start in Section 2 with a description of the organisation and structure of the shared task, along with track details. Next, we summarise results at the level of individual experiments, in terms of the reproduction task, and different degree-of-reproducibility assessments, first for Track B (Section 3), then Track A (Section 4).

In Section 5, we look at the quality criteria assessed in evaluations and other properties of the ReproNLP evaluation studies in standardised terms as facilitated by HEDS datasheets, and explore if any of these show signs of affecting degree of reproducibility (Section 5). We conclude with some discussion (Section 6) and a look to future work (Section 7).

## 2. ReproNLP 2024

ReproNLP 2024[1] consisted of two tracks, one an 'unshared task' in which teams re-run their own or any other previous work (Track A), the other a standard shared task in which teams re-run one of a set of organiser-selected experiments (Track B):

A **Open Track:** Repeat any previously reported work developing and evaluating systems, and report the approach and outcomes. Unshared task.

B **ReproHum Track**: For a shared set of selected evaluation studies (listed below) from the ReproHum Project, participants repeat one or more of the studies and compare results, using the information provided by the ReproNLP organisers only, and following a common reproduction approach.

Track B forms part of the ReproHum project[2] and the studies offered in it were selected according to criteria of suitability and balance to form part of a larger coordinated multi-lab multi-test reproduction study, as described in detail elsewhere (Belz et al., 2023).

---

[1] All information and resources relating to ReproNLP are available at https://repronlp.github.io/.
[2] https://reprohum.github.io/

An overview of the papers we selected experiments from, and the complete studies the latter formed part of, is presented below. Note that we only include here the original papers for which we received submissions; there were 21 papers offered in the track in total (the full list can be found on the ReproNLP website).

The information provided for each study below covers whether the assessment of systems was *relative* to other systems or *absolute* without comparitors; what the language*(s)* of the systems were; how many *datasets* were used; how many *systems* were evaluated and by how many *evaluators*; and whether the evaluation was run on a *crowdsourcing* platform.

1. **Reif et al. (2022)**: *A Recipe for Arbitrary Text Style Transfer with Large Language Models*: https://aclanthology.org/2022.acl-short.94

   Absolute evaluation study; English; 3 quality criteria; 3 datasets; varies between 4 and 6 systems and between 200 and 300 evaluation items per dataset-criterion combination; crowdsourced.

2. **Liu et al. (2021)**: *DExperts: Decoding-Time Controlled Text Generation with Experts and Anti-Experts*: https://aclanthology.org/2021.acl-long.522

   Relative evaluation study; English; 3 quality criteria; 2 datasets; varies between 5 and 6 systems and between 960 and 1200 evaluation items per dataset-criterion combination; crowdsourced.

3. **Atanasova et al. (2020)**: *Generating Fact Checking Explanations*: https://aclanthology.org/2020.acl-main.656

   Absolute evaluation study; English; 1 quality criterion; 1 dataset; 3 systems and 240 evaluation items. Relative evaluation study; English; 4 quality criteria; 1 dataset; 3 systems and 40 evaluation items per criterion.

4. **August et al. (2022)**: *Generating Scientific Definitions with Controllable Complexity*: https://aclanthology.org/2022.acl-long.569

   Absolute evaluation study; English; 5 quality criteria; 2 datasets; 3 systems and 300 evaluation items per dataset-criterion combination; some crowdsourced.

5. **Hosking et al. (2022)**: *Hierarchical Sketch Induction for Paraphrase Generation*: https://aclanthology.org/2022.acl-long.178

   Relative evaluation study; English; 3 quality criteria; 1 dataset; 4 systems and 1800 evaluation items per criterion; crowdsourced.

6. **Yao et al. (2022)**: *It is AI's Turn to Ask Humans a Question: Question-Answer Pair Generation for Children's Story Books*: https://aclanthology.org/2022.acl-long.54

   Absolute evaluation study; English; 3 quality criteria; 1 dataset; 3 systems and 361 evaluation items per criterion.

7. **Feng et al. (2021)**: *Language Model as an Annotator: Exploring DialoGPT for Dialogue Summarization*: https://aclanthology.org/2021.acl-long.117

   Absolute evaluation study; English; 3 quality criteria; 2 datasets; 7 systems and varies between 70 and 700 evaluation items per dataset-criterion combination.

8. **Gabriel et al. (2022)**: *Misinfo Reaction Frames: Reasoning about Readers' Reactions to News Headlines*: https://aclanthology.org/2022.acl-long.222

   Absolute evaluation study; English; 3 quality criteria; 1 dataset; 3 systems and 588 evaluation items per criterion; crowdsourced.

9. **Kasner & Dusek (2022)**: *Neural Pipeline for Zero-Shot Data-to-Text Generation*: https://aclanthology.org/2022.acl-long.271

   Absolute evaluation study; English; 5 quality criteria; 2 datasets; 6 systems and 600 evaluation items per dataset-criterion combination.

10. **Shardlow & Nawaz (2019)**: *Neural Text Simplification of Clinical Letters with a Domain Specific Phrase Table*: https://aclanthology.org/P19-1037

    Relative evaluation study; English; 1 quality criterion; 1 dataset; 4 systems and 100 evaluation items; crowdsourced.

11. **Castro Ferreira et al. (2018)**: *NeuralREG: An end-to-end approach to referring expression generation*: https://aclanthology.org/P18-1182

    Absolute evaluation study; English; 3 quality criteria; 1 dataset; 6 systems and 144 evaluation items per criterion; crowdsourced.

In the ReproHum multi-lab multi-test study (for which the above papers were selected), rather than attempt to repeat entire studies, we decided to use our limited resources to repeat assessments of individual quality criteria on individual datasets (which is what we mean by a single 'experiment'), with specific properties so as to have equal numbers of assessments with the specific properties the ReproHum study is designed to compare. Some of the properties of these individual experiments are given in Table 2 alongside the (single) quality criteria they assess.

Each of these experiments is being re-run in two separate reproduction studies in ReproHum. Those that have completed in the current batch are

included here in the ReproNLP'24 report. All experiments from the current and preceding batch (the latter reported in ReproNLP'23) were also open to all other ReproNLP'24 participants.

Note that non-ReproHum participants were free to include more than the ReproHum experiment in their reproduction study, and some did (Section 4).

We obtained agreement from the original authors to use their experiments in the ReproHum project and provided very detailed information about the experiments which were shared with all participants.

## 2.1. Participation

There were three submissions for Track A and 15 for Track B. One submission in Track A did not meet our quality threshold and was rejected. The ReproHum partners reporting in Track B are listed in Table 1. The non-ReproHum participating labs were University of Bucharest (Florescu et al., 2024) in Track B, and Heriot-Watt University (Sasidharan Nair et al., 2024) and ADAPT Centre / Dublin City University (Lorandi and Belz, 2024) in Track A (see Sections 4 and 3.2, respectively).

## 2.2. Approach to reproduction and reproducibility assessment

We encouraged all participants to complete a HEDS datasheet (Shimorina and Belz, 2022) in the ReproHum version,[3] and to follow the ReproHum Common Approach to reproduction laid out in Appendix A which includes QRA++ (Belz, 2022; Belz and Thomson, 2023), an approach to measuring how close results from two evaluations are, and how reproducible evaluation measures are, in a way that accommodates multiple reproduction studies of the same original work and is comparable across different such sets of reproductions.

In this report we analyse all submissions in terms of QRA++ measures recomputed by us to facilitate comparison across submissions. In brief summary, QRA++ distinguishes four types of results commonly reported in NLP and ML papers:

1. Type I results: single numerical scores, e.g. mean quality rating, error count, etc.

2. Type II results: sets of related numerical scores, e.g. set of Type I results .

3. Type III results: categorical labels attached to text spans of any length.

4. Type IV results: Qualitative findings stated explicitly or implied by quantitative results in the original paper.

The above are quantitatively assessed as follows:

1. Type I results: Small-sample coefficient of variation CV* (Belz, 2022).

2. Type II results: Pearson's r, Spearman's $\rho$.

3. Type III results: Multi-rater: Fleiss's $\kappa$; Multi-rater, multi-label: Krippendorff's $\alpha$.

4. Type IV results: Proportion of findings that are / are not confirmed by the repeat experiment. To obtain comparable results we restrict ourselves to pairwise system ranks as findings.

In the submissions analysed in this paper we have Type I, II and IV results, and therefore apply the corresponding quantitative measures above. CV* plays a central role in our analyses, and is a version of the standard coefficient of variation corrected for small samples (Belz, 2022).

The ReproHum reproduction studies were strictly controlled to be comparable to each other and the original work. However, there was a difference between the studies reported in 2023 and 2024 in this respect. For the earlier batch, our aim was to achieve maximum similarity between original and reproduction studies, and we strove to resolve every last bit of lack of clarity. In the batch reported here, we abandoned this ultimately infeasible approach, recognising that evaluation experiments should be robust to minor differences. As a result, when there was insufficient clarity about how an aspect of an experiment was implemented, partner labs drafted solutions which were moderated by the ReproHum project team to provide an agreed solution that both partner labs reproducing the same experiment then used. For more details on such cases, please see the individual submission reports in this volume.

Finally, we have by now gathered a sufficient number of reproduction studies reporting CV* values to support the following categorisation for *human evaluations*: we refer to any CV* from 0 to around 10 as indicating a good degree of reproducibility, between 10 and around 30 as medium, and anything above that as poor.

Note that high CV* scores indicate poor reproducibility, and vice versa.

## 3. Track B

The subsections within Sections 3.1 and 3.2 each report the results from all reproduction studies for one of the Track B experiments, Sections 3.1 as conducted by ReproHum partners, in Sections 3.2 by other ReproNLP participants.

In each such subsection, we start by giving a brief summary of the experiment. Next, we show the system-level evaluation scores from the original study and the either one or two reproduction studies, alongside the corresponding CV* value computed on all either two or three scores. We finish

---

[3] https://github.com/nlp-heds/repronlp2024

| Original Study | Qual. Criterion | #ev-ors | #sys | items-per-sys | Labs reproducing study |
|---|---|---|---|---|---|
| Liu et al. (2021) | Fluency | varies | 5 | 192 | a) Heriot-Watt University <br> b) U. de Santiago de Compostela |
| Hosking et al. (2022) | Preservation of meaning | varies | 4 | 450 | a) University of Illinois Chicago <br> b) Edinburgh Napier University |
| Feng et al. (2021) | Informativeness | 4 | 7 | 10 | a) Bielefeld University <br> b) Charles University |
| Atanasova et al. (2020) | Coverage | 3 | 3 | 13 | a) Manchester University <br> b) Peking University |
| August et al. (2022) | Fluency | 2 | 3 | 100 | a) Tilburg University <br> b) University of Groningen |
| Castro Ferreira et al. (2018) | Clarity | 60 | 6 | 24 | a) trivago |
| Kasner and Dusek (2022) | Number of redundancies per system | 2 | 6 | 100 | a) Technological University Dublin |
| Shardlow and Nawaz (2019) | Ease of understanding | varies | 4 | 25 | a) University of Groningen |
| Gabriel et al. (2022) | Social acceptability | varies | 3 | 196 | a) University of Cape Town |

Table 1: ReproNLP experiments performed by ReproHum partner labs. All experiments were in the English language. The number of evaluators sometimes varies because some original studies did not control for this property, but rather allowed as many crowd-source participants to rate as many items as they wished. An item is defined as one system output evaluated absolutely, or a set of system outputs evaluated relatively.

by reporting the pairwise Pearson's $r$ and Spearman's $\rho$ correlation coefficients (Type II QRA) and the proportion of findings upheld (Type IV QRA). (See also Section 2.2.) In the present context, we consider each pairwise system ranking to be one finding. All scores are recomputed by us from the results reported in participants' papers, and those in the original studies.

As noted above, we report Type I, II, and IV QRA results only. This is because in most cases there are no Type III results, and in some cases where there are Type III results we do not have access to all of the raw annotations from the original studies (which would be needed in order to calculate Type III QRA).

## 3.1. Track B: ReproHum Partners

In this section, we summarise results from the reproduction studies performed by ReproHum partner labs reporting in Track B. We have five pairs of such studies, and four single studies where a second lab has either not yet completed and/or been assigned.

### 3.1.1. Liu et al. (2021)

In this experiment, participants were shown pairs of outputs from a new controlled text generation system (DExperts) and four different baselines. They were then asked which is more **fluent**. The follow-

ing table shows the proportion of times DExperts was preferred over (>), considered equally good as (=), or dispreferred (<), over each of the four baselines, in original study (abbreviated O), reproduction 1 or *R1* (Dinkar et al., 2024), and reproduction 2 or *R2* (González Corbelle et al., 2024). The highest such proportion is highlighted in boldface. The last column shows the corresponding CV* (n=3) values for each row, finding overall a medium to poor degree of reproducibility.

| System | O | R1 | R2 | CV* |
|---|---|---|---|---|
| DExperts > GPT-2 | 0.30 | **0.39** | **0.35** | 15.90 |
| GPT-2 = DExperts | **0.40** | 0.23 | 0.32 | 32.83 |
| GPT-2 > DExperts | 0.30 | 0.38 | 0.33 | 14.67 |
| DExperts > DAPT | 0.26 | **0.42** | 0.30 | 31.16 |
| DAPT = DExperts | **0.39** | 0.19 | 0.29 | 42.15 |
| DAPT > DExperts | 0.35 | 0.40 | **0.41** | 10.16 |
| DExperts > PPLM | **0.37** | **0.47** | **0.39** | 15.78 |
| PPLM = DExperts | 0.33 | 0.19 | 0.28 | 32.52 |
| PPLM > DExperts | 0.31 | 0.33 | 0.33 | 4.37 |
| DExperts > GeDi | **0.36** | **0.45** | **0.36** | 16.29 |
| GeDi = DExperts | 0.35 | 0.20 | 0.29 | 32.96 |
| GeDi > DExperts | 0.28 | 0.35 | 0.35 | 15.12 |
| Mean CV* | – | – | – | 21.99 |

In terms of Type II QRA, the correlations between each pair of columns above are as shown in the next table below. We can see that both $r$ and $\rho$ are negative for O and R1, and around 0 for O and R2.

In stark contrast, they are both medium to strong in the positive direction for R1 and R2.

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---------|---------|-------|-------|---------|
| O | R1 | -0.36 | -0.18 | 2/4 |
| O | R2 | 0.07 | 0.01 | 2/4 |
| R1 | R2 | 0.75 | 0.79 | 3/4 |

The above table also includes Type IV assessment, which assesses the proportion of times the pairwise system rank (e.g. DExperts was found to be better than PPLM) was upheld by a reproduction experiment. For this particular experiment we determined pairwise system rank as the relationship (>, <, =) that was selected most often by participants for a given pair of systems. In this way, we can see that both reproductions confirmed 50% (2/4) findings from the original experiment (the same two in both cases), while they agreed more with each other than the original study.

### 3.1.2. Hosking et al. (2022)

Here, participants were shown pairs of outputs from paraphrase generation systems and asked which best **preserves the meaning** of the input text. The below table shows scores that represent the strength with which a system was (dis)preferred on a scale from -100 to +100 (negative meaning dispreferred), alongside the corresponding CV* (n=3) values, for *O* (the original study), *R1* (Arvan and Parde, 2024), and *R2* (Watson and Gkatzia, 2024), finding a good degree of reproducibility at the level of system scores, with uniformly low CV*.

| System | O | R1 | R2 | CV* |
|--------|-------|--------|--------|------|
| VAE | 36.00 | 37.04 | 23.00 | 7.24 |
| Latent BoW | -16.00 | -14.52 | -8.67 | 5.45 |
| Separator | -24.00 | -29.78 | -17.89 | 9.55 |
| HRQ-VAE | 4.00 | 7.26 | 3.56 | 2.35 |
| Mean CV* | – | – | – | 6.15 |

The correlations (Type II QRA) between all experiments are near perfect, and the pairwise ranks of systems (Type IV QRA) are confirmed in all cases:

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---------|---------|------|------|---------|
| O | R1 | 0.99 | 1.00 | 6/6 |
| O | R2 | 0.99 | 1.00 | 6/6 |
| R1 | R2 | 0.99 | 1.00 | 6/6 |

With all three QRA measures across both reproductions strongly confirming the original results, Hosking et al. has one of the three highest overall degree of reproducibility of any of the human evaluations in ReproNLP'24 (the other two being for Shardlow & Narwaz for ease of understanding, and Yao et al. for readability, below).

### 3.1.3. Feng et al. (2021)

For this experiment, participants were asked to rate system outputs on a scale of 1 (worst) to 5 (best) the **informativeness** of paragraph-sized summaries of multi-page meeting transcriptions. The below table shows the mean system scores from *O* (the original study), *R1* (Fresen et al., 2024), and *R2* (Lango et al., 2024), alongside the corresponding CV* (n=3) values, showing reproducibility for system scores across the board.

| System | O | R1 | R2 | CV* |
|--------|------|------|------|-------|
| Golden | 4.70 | 2.40 | 4.60 | 54.80 |
| PGN | 2.92 | 2.18 | 1.53 | 70.26 |
| HMNet | 3.52 | 2.20 | 2.68 | 45.37 |
| PGN(DKE) | 3.20 | 2.18 | 1.93 | 57.24 |
| PGN(DRD) | 3.15 | 3.00 | 1.90 | 49.56 |
| PGN(DTS) | 3.05 | 2.27 | 1.85 | 53.55 |
| PGN(DALL) | 3.33 | 2.52 | 1.85 | 57.83 |
| Mean CV* | – | – | – | 55.52 |

The next table below shows that despite much lower scores for all but the 'Golden' system, strong correlations are seen between the original study (O) and R2. However, correlations between O and R1, and between R1 and R2, are close to 0 (no correlation).

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---------|---------|-------|------|---------|
| O | R1 | 0.01 | 0.27 | 12/21 |
| O | R2 | 0.99 | 0.85 | 18/21 |
| R1 | R2 | -0.03 | 0.11 | 11/21 |

This picture is somewhat confirmed by the Type IV QRA scores which show best confirmation of results for R2 but interestingly also show that R1, despite the other QRA results above, still confirmed about half the findings from O.

### 3.1.4. Atanasova et al. (2020)

Here, participants were asked to rank the justifications generated by three different systems in terms of their **coverage** relative to an input claim. The below table shows the mean rank for each system from *O* (the original study), *R1* (Loakman and Lin, 2024), and *R2* (Gao et al., 2024), alongside the corresponding CV* (n=3) values. The latter show the degree of reproducibility of the mean system rank to be good to medium for the two Explain systems, but poor for the Just system.[4]

---

[4]Note that here we have a question mark over the accuracy of the scores reported in the original study. We had the raw responses from the original experiment available to us and both reproducing teams recalculated system scores on this basis, with neither team matching the original results (or each other).

| System | O | R1 | R2 | CV* |
|---|---|---|---|---|
| Just | 1.48 | 1.62 | 2.18 | 59.58 |
| Explain-Extr | 1.89 | 2.05 | 1.93 | 10.64 |
| Explain-MT | 1.68 | 1.78 | 1.62 | 14.25 |
| Mean CV* | – | – | – | 28.16 |

As the Type II/IV table below shows, strong correlations were found, and all findings were confirmed, between *O* (the original study) and *R1*. However, both QRA measures were very poor for O and R2, and also R1 and R2.

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.99 | 1.00 | 3/3 |
| O | R2 | -0.43 | -0.50 | 1/3 |
| R1 | R2 | -0.31 | -0.50 | 1/3 |

### 3.1.5. August et al. (2022)

For this experiment, participants were asked to rate the **fluency** of generated scientific definitions on a scale of 1 (not at all) to 4 (very). The below table shows the mean system scores, alongside the corresponding CV* (n=3) values, for *O* (the original study), *R1* (van Miltenburg et al., 2024), and *R2* (Li et al., 2024), finding a medium to borderline poor degree of reproducibility for all systems, albeit better for the SVM system.

| System | O | R1 | R2 | CV* |
|---|---|---|---|---|
| SVM | 3.71 | 3.12 | 3.02 | 19.96 |
| GeDi | 3.20 | 2.57 | 2.40 | 29.90 |
| DExpert | 2.33 | 2.28 | 1.81 | 30.76 |
| Mean CV* | – | – | – | 26.87 |

Correlations were very strong between all studies, with the order of pairwise ranks (Type IV) confirmed in all cases:

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.95 | 1.00 | 3/3 |
| O | R2 | 0.99 | 1.00 | 3/3 |
| R1 | R2 | 0.99 | 1.00 | 3/3 |

### 3.1.6. Castro Ferreira et al. (2018)

Participants were shown outputs of a data-to-text system and asked to rate their **clarity** on a 1 (very bad) to 7 (very good) scale. The below table shows the mean system ratings, alongside the corresponding CV* (n=2) values, for *O* (the original study) and *R1* (Mahamood, 2024), finding an excellent degree of reproducibility across the board.

| System | O | R1 | CV* |
|---|---|---|---|
| OnlyNames | 4.90 | 4.92 | 0.51 |
| Ferreira | 4.93 | 4.69 | 6.28 |
| NeuralREG+Seq2Seq | 4.97 | 4.97 | 0.00 |
| NeuralREG+CAtt | 5.26 | 4.97 | 7.03 |
| NeuralREG+HierAtt | 5.13 | 5.04 | 2.20 |
| Original | 5.42 | 5.22 | 4.62 |
| Mean CV* | – | – | 3.44 |

Correlations between R1 and O were medium-strong, and 80% (12/15) of pairwise system rankings were confirmed:

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.78 | 0.84 | 12/15 |

### 3.1.7. Kasner and Dusek (2022)

This experiment was originally an error analysis performed by the authors, although it fits the definition of a human evaluation used in the ReproHum Project. Participants (the two authors) were shown the input and outputs from data-to-text systems and asked to count the **number of repetitions** in the outputs. The below table shows repetition error counts for different systems and corresponding CV* (n=2) values for *O* (the original study) and *R1* (Klubička and Kelleher, 2024), finding extremely poor degrees of reproducibility at the level of system scores.

| System | O | R1 | CV* |
|---|---|---|---|
| Full-3-Stage | 0 | 13 | 199.40 |
| Full-2-Stage | 1 | 11 | 166.17 |
| Full-1-Stage | 79 | 156 | 65.34 |
| Filtered-3-Stage | 0 | 9 | 199.40 |
| Filtered-2-Stage | 0 | 10 | 199.40 |
| Filtered-1-Stage | 41 | 84 | 68.59 |
| Mean CV* | – | – | 149.72 |

However, Pearson's is very nearly perfect,[5] with Spearman's a less strong 0.82, and 73% of pairwise system ranks confirmed:

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.99 | 0.82 | 11/15 |

We thus have a mixed picture here with system score level reproducibility extremely poor, about three quarters of findings confirmed, a reasonably strong rank correlation and near perfect product-moment correlation.

---

[5]It would round up to 1.00 but our rounding policy keeps it at 0.99 to avoid giving a false impression. See Appendix B.

### 3.1.8. Shardlow and Nawaz (2019)

In this experiment, participants were shown medical texts, from four text simplification systems, and asked to rank them from best to worst in terms of **ease of understanding**. The below table shows mean system rank for each system and the corresponding CV* (n=2) values for *O* (the original study) and *R1* (Li et al., 2024), finding a good degree of reproducibility for mean rank when considering the two *NTS* systems, and an excellent degree for the *ORIG* and *PTB* systems.

| System | O | R1 | CV* |
|---|---|---|---|
| NTS+PT | 1.93 | 1.82 | 12.53 |
| NTS | 2.34 | 2.46 | 8.55 |
| ORIG | 2.79 | 2.76 | 1.69 |
| PTB | 2.94 | 2.96 | 1.02 |
| Mean CV* | – | – | 5.95 |

There were also strong correlations between the two studies and the pairwise ranks were confirmed in all cases:

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.98 | 1.00 | 6/6 |

With all three QRA measures across both reproductions strongly confirming the original results, Shardlow & Narwaz has one of the three highest overall degree of reproducibility of any of the human evaluations in ReproNLP'24 (the other two being for Hosking et al. for meaning preservation, above, and Yao et al. for readability, below).

### 3.1.9. Gabriel et al. (2022)

For this experiment, participants were shown output texts from three systems that generate statements of the writer's intents given news headlines as input. Their task was to decide whether the text was socially acceptable or not. The below table shows the percentage of times a system was deemed socially acceptable alongside the corresponding CV* (n=2) values for *O* (the original study) and *R1* (Mahlaza et al., 2024), finding a good to medium degree of reproducibility.

| System | O | R1 | CV* |
|---|---|---|---|
| T5-base | 75.30 | 68.67 | 9.18 |
| T5-large | 74.66 | 68.31 | 8.86 |
| GPT-2 (large) | 74.66 | 65.30 | 13.34 |
| Mean CV* | – | – | 10.46 |

Pearson's $r$ was only moderate, with a stronger Spearman's $\rho$, and 67% of findings confirmed:

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.58 | 0.87 | 2/3 |

### 3.2. Track B: Other teams

Track B of ReproNLP was also open to non-ReproHum partner labs. Participants in this track reproduce experiments of their choice from the same set of Track B papers, but do not necessarily follow the exact common approach (Appendix A) as ReproHum partner labs do.

### 3.2.1. Yao et al. (2022)

Florescu et al. (2024) repeated this evaluation of generated questions and answers for children's stories, performing the evaluation for all three quality criteria in the original study.

**Readability:** Participants are asked to rate what was named the "readability" of the question-answer pair. The exact prompt used, however, was "readability(grammarly [sic] correct and clear language. worst 1 to 5)", which references three different quality criteria (readability, grammaticality and clarity), making it a clear example of the confusion in quality criteria found by Howcroft et al. (2020).

The below table shows the mean readability ratings for each system alongside CV* (n=2) values for the original study *(O)* and *R1* (Florescu et al., 2024), finding a good degree of reproducibility.

| System | O | R1 | CV* |
|---|---|---|---|
| Ours | 4.71 | 4.52 | 5.24 |
| PAQ Baseline | 4.08 | 4.17 | 2.87 |
| Groundtruth | 4.95 | 4.71 | 6.25 |
| Mean CV* | – | – | 4.79 |

Correlations were near perfect, and the pairwise ranks of systems were confirmed in all cases:

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.99 | 1.00 | 3/3 |

With all three QRA measures strongly confirming the original results, Yao et al. has one of the three highest overall degree of reproducibility of any of the human evaluations in ReproNLP'24 (the other two being for Hosking et al. for meaning preservation, and Shardlow & Narwaz for ease of understanding, above).

**Relevancy (Question):** The following table shows the mean question relevancy ratings for systems alongside the corresponding CV* (n=2) values for the original study *(O)* and *R1* (Florescu et al., 2024), finding only a medium degree of reproducibility for the system outputs, with a good degree of reproducibility for the human-authored ground truth.

| System | O | R1 | CV* |
|---|---|---|---|
| Ours | 4.39 | 3.83 | 17.95 |
| PAQ Baseline | 4.18 | 3.61 | 19.63 |
| Groundtruth | 4.92 | 4.71 | 5.49 |
| Mean CV* | – | – | 14.36 |

However, correlations are still perfect, with the pairwise ranks of systems confirmed in all cases:

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.99 | 1.00 | 3/3 |

**Relevancy (Answer):**  Finally, for the relevancy of the answer, system scores are again less reproducible than the ground truth, with a medium to poor degree of reproducibility for the systems, and good degree of reproducibility for the ground truth:

| System | O | R1 | CV* |
|---|---|---|---|
| Ours | 3.99 | 3.20 | 30.35 |
| PAQ Baseline | 3.90 | 3.20 | 27.37 |
| Groundtruth | 4.83 | 4.46 | 10.12 |
| Mean CV* | – | – | 22.61 |

Correlations are strong with Spearman's lower, and the two systems being scored identically in the reproduction experiment, as opposed to only being similar in the original study; this also affects the pairwise rankings confirmed score (Type IV) for the two systems:

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.99 | 0.87 | 2/3 |

# 4.  Track A

We accepted two submissions in the open track, where participants could carry out reproduction experiments for any paper, focusing on human and/or metric-based evaluations.

## 4.1.  Chakravarthi et al. (2020)

In the original study, a code-mixed Malayalam language dataset was annotated for sentiment (5 labels) by human participants and then used to train classifiers which were in turn evaluated by automatic metrics. Sasidharan Nair et al. (2024) recreate this complete pipeline.

### 4.1.1.  Label counts (human evaluation)

The count of labels recorded in the reproduction varied greatly from the original study, resulting in a moderate degree of reproducibility for some labels, and a very poor degree of reproducibility for others, as shown in this table:

| System | O | R1 | CV* |
|---|---|---|---|
| Positive | 565 | 626 | 10.21 |
| Negative | 138 | 162 | 15.95 |
| Mixed Feelings | 70 | 144 | 68.95 |
| Neutral | 398 | 327 | 19.53 |
| Non-Malayalam | 177 | 89 | 65.97 |
| Mean CV* | – | – | 36.12 |

However, the correlations were strong, with pairwise ranks also confirmed for most labels. Note that rather than comparing systems we are comparing label counts of an annotated corpus.

| Study A | Study B | $r$ | $\rho$ | Type IV |
|---|---|---|---|---|
| O | R1 | 0.94 | 0.70 | 8/10 |

### 4.1.2.  Automated metrics

After completing their re-annotation of the corpus, Sasidharan Nair et al. (2024) then evaluated LR and BERT sentiment classifiers on both the original corpus and their newly created one, using F1 score. The below table shows Mean CV* for *O* (results from the original paper) and *Re-Imp* (a re-implemented classifier by Sasidharan Nair et al. (2024) but trained on the original corpus). The table also shows the Mean CV* for *Re-Imp* and *Re-Ann*, where *Re-Ann* refers to the re-implemented classifier trained on the re-annotated corpus.

This reproduction study clearly shows the effect that the reproducibility of human evaluation can have on the reproducibility of downstream tasks.

| Classifier | Study A | Study B | Mean CV* |
|---|---|---|---|
| LR | O | Re-imp | 7.65 |
| BERT | O | Re-imp | 22.73 |
| LR | Re-imp | Re-ann | 47.70 |
| BERT | Re-imp | Re-ann | 24.10 |

Note that we calculate Mean CV* in this report differs from how Sasidharan Nair et al. (2024) calculate it, where the macro and weighted averages of F1 score are calculated first, with the mean CV* then calculated at that level.

In terms of Type IV results, the reproducing team find that by both macro-average and weighted-average, for both setups (*O* vs *Re-imp* and *Re-Imp* vs *Re-ann*), the BERT classifier is always better than LR. This corresponds to 8/8 findings upheld. Even at the per-label level, 9/10 findings are upheld for *O* vs *Re-imp*, and 8/10 for *Re-Imp* vs *Re-ann*.

## 4.2.  Gu et al. (2022, 2023)

Lorandi and Belz (2024) reproduce original studies found in Gu et al. (2022) and Gu et al. (2023). They calculate the CV* between original and reproduc-

| | ReproNLP 2024 | | | | | | mean CV* | | |
|---|---|---|---|---|---|---|---|---|---|
| Orig Study // *Repro a / Repro b*<br>**measurands** | **3.2.1** | **4.3.4** | **4.3.8** | **4.1.1** | **4.1.2** | **4.1.3** | **a(n=2)** | **b(n=2)** | **n=3** |
| Liu et al. (2021) // *Dinkar et al. (2024) / González Corbelle et al. (2024)*<br>**Fluency** | UNK / 96 / 90 | A,B,Tie | RQE | Good | Both | iiOR | 34.55 | 14.58 | 21.99 |
| Hosking et al. (2022) // *Arvan and Parde (2024) / Watson and Gkatzia (2024)*<br>**Preservation of meaning** | UNK / 180 / 180 | A,B | RQE | Good | Cont | RtI | 3.37 | 6.62 | 6.15 |
| Feng et al. (2021) // *Fresen et al. (2024) / Lango et al. (2024)*<br>**Informativeness** | 4 / 4 / 4 | 1-5 | DQR | Good | Cont | RtI | 52.07 | 70.53 | 55.52 |
| Atanasova et al. (2020) // *Loakman and Lin (2024) / Gao et al. (2024)*<br>**Coverage** | 3 / 3 / 3 | 1-3 | RQE | Good | Cont | RtI | 18.49 | 32.56 | 28.16 |
| August et al. (2022) // *van Miltenburg et al. (2024) / Li et al. (2024)*<br>**Fluency** | 2 / 2 / 2 | 1-4 | DQE | Good | Both | iiOR | 20.50 | 40.62 | 26.87 |
| Castro Ferreira et al. (2018) // *Mahamood (2024)*<br>**Clarity** | 60 / 60 | 1-7 | DQE | Good | Both | iiOR | 3.44 | - | - |
| Kasner and Dusek (2022) // *Klubička and Kelleher (2024)*<br>**Number of redundancies per system** | 2 / 2 | count | Count | Good | Cont | iiOR | 149.72 | - | - |
| Shardlow and Nawaz (2019) // *Mondella et al. (2024)*<br>**Ease of understanding** | 98 / 40 | 1-4 | RQE | Good | Both | iiOR | 5.95 | - | - |
| Gabriel et al. (2022) // *Mahlaza et al. (2024)*<br>**Social acceptability** | UNK / 42 | Yes,No | DQE | Feature | Both | EFoR | 10.46 | - | - |

Table 2: Summary of some properties of ReproNLP experiments performed by ReproHum partner labs, alongside mean CV* (n=2, or n=3; shown in different columns because different samples sizes are not directly comparable). 3.2.1 = number of evaluators in original/reproduction experiment; 4.3.4 = List/range of possible responses; 4.3.8 = Form of response elicitation (DQE: direct quality estimation, RQE: relative quality estimation, Cl/Lab: classification/labelling, Count: counting occurrences in text); 4.1.1 = Correctness/Goodness/Features; 4.1.2 = Form/Content/Both; 4.1.3 = each output assessed in its own right (iiOR) / relative to inputs (RtI) / relative to external reference (EFoR).

tion scores for each evaluation measure. Based on these, we include per-system mean CV* scores below, along with the maximum and minimum.

| System | CV* | | |
|---|---|---|---|
| | Mean | Min | Max |
| Multi-CTG | 1.68 | 0.34 | 5.56 |
| Prior-CTG | 1.25 | 0.00 | 4.14 |
| Prior-CTG+optim | 1.28 | 0.00 | 3.81 |

Moreover, Lorandi and Belz (2024) report correla-

tions (Type II) in excess of 0.99 for all reproductions and all findings (Type IV) as upheld.

A full breakdown of per-measure CV* results, as well as other analyses, can be found in their paper.

## 5. Reproducibility by Quality Criterion and other properties

We saw a wide variety of different degrees of reproducibility for the different human evaluations

in previous sections. It seems likely that these differences in degree of reproducibility are explainable by differences between the evaluations. In the QRA++ approach (Belz, 2022; Belz and Thomson, 2023), as in metrology on which it is based, such differences are captured by 'conditions of measurement,' and HEDS was designed to capture these.

Table 2 shows some of the main HEDS properties of the experiments repeated by ReproHum partner labs, along with mean CV* values calculated as follows:

- **a(n=2)**: the mean of two-way CV* values between *Orig Study* and *Repro a*.

- **b(n=2)**: the mean of two-way CV* values between *Orig Study* and *Repro b* (if there was a *Repro b*).

- **(n=3)**: the mean of three-way CV* values between *Orig Study*, *Repro a*, and *Repro b* (if there were 3 sets of results).

What we are looking for in this table is any indication that one of the HEDS properties affects experiment-level mean CV* (last three columns). One such property is number of evaluators (HEDS Question 3.2.1): the pattern is for larger number of evaluators to be associated with better reproducibility, with the exception of Fluency in Liu et al. which bucks the trend somewhat.

Another trend that is observable is that evaluations that are more cognitively complex tend to have poorer reproducibility than cognitively simpler evaluations. An extreme example of this is Kasner and Dusek's count of redundancies per system, which are very hard to match in reproductions. Similar results were obtained in an earlier pair of reproductions of an error analysis experiment, where some of the error counts also reached above 140 CV* (**?**). Another example is Informativeness (fourth from top in table).

Cognitively simpler assessments like Clarity and Fluency have better score-level reproduciblity. This is a trend that we have consistently observed across multiple reproduction experiments. Note however that here too Fluency in Liu et al. bucks the trend which may be explained by other experimental properties we are not examining here.

## 6. Discussion

As we saw in previous sections, different types of QRA++ assessments (Type I, II, and IV) can show very different degrees of reproducibility for sets of reproductions for the *same* original experiment. For example, for Social Acceptability in Gabriel et al. 2022, CV* levels were reasonable but Pearson's was only 0.58.

Another example is Fluency in August et al. (2022) where the CV* values are quite poor, but Type II and IV reproducibility is excellent.

It can also be the case that one reproduction for the same original experiment indicates excellent reproducibility and another shows very poor reproducibility, as was the case for Content Coverage Atanasova et al. 2020.

The latter observation (observed previously) indicates the importance of conducting more than one reproduction experiment. An alternative may be to increase the number of individual assessments carried out (Simonsohn, 2015), but it is not clear how additional assessments should be created (more evaluators, more system outputs, both?).

The differences between results from different types of QRA assessment highlight that each assesses a different aspect of reproducibility: Type I/CV* looks at how close individual aligned scores are; Type II/correlations look at how similar relative increases and decreases are in aligned sets of scores; and Type IV/findings abstracts away from scores altogether to look at findings which we here interpret as pairwise system ranks, i.e. which of two systems performs better.

Ultimately, it is the latter, pairwise system ranks, that we care most about in many contexts. What matters is not necessarily maximising the rank correlation, but the proportion of pairwise ranks that are the same (although clearly these are linked).

In the previous section we looked at the effect different experiment properties may have on reproducibility. However, these cannot explain differences between reproductions of the same original experiment where properties are the same. This means that there are other factors affecting reproducibility, e.g. evaluator sampling, and quality of the reproduction experiment. All of this clearly makes it harder to link properties with reproducibility.

Given the finding that score-level QRA can show poor reproducibility even where all findings are upheld, it might be questioned whether interannotator agreement (IAA), commonly used as an indicator of experiment quality, is really the right measure. It might be that reproducibility tests assessing multiple different QRA measures are a better pre-experiment test of quality.

## 7. Conclusion

Shared task result reports tend to be written under considerable pressure of time, and the present paper is no exception. We will conduct additional analyses and more in-depth explorations of our data in due course, as well as reporting the results from the second batch of ReproHum multi-lab multi-test study experiments once all have been completed. The latter will provide more robustly

quantifiable assessment of the impact of selected experiment properties on reproducibility.

This year's edition of the shared task has once again highlighted the considerable extent to which results (i) from different reproduction experiments of the same original experiments, and (ii) from different types of QRA analysis, can differ. This can be interpreted as meaning that we should conduct multiple reproduction experiments, and multiple types of QRA analysis, respectively.

There continues to be little standardisation in evaluation practices, and quality criteria names and definitions in particular, in human evaluation in NLP, despite numerous surveys and studies (Belz et al., 2020; Howcroft et al., 2020; van der Lee et al., 2021; Gehrmann et al., 2023) calling for more standardisation to improve quality and reliability. In the present context, lack of standardisation also has the effect of muddying the waters with respect to conclusions about which quality criteria are associated with better reproducibility: if it is unclear, e.g. due to mere name differences, whether the same quality criterion was in fact assessed, it is hard to draw accurate conclusions beyond the individual experiment.

All in all, it seems clear that human evaluation in NLP would benefit from more standardisation in experimental design and execution, for better comparability, but also so that reproducibility, hence reliability, of standard methods can be established, and once established, benefited from.

## Acknowledgments

## Bibliographical References

Mohammad Arvan and Natalie Parde. 2024. Human evaluation reproduction report for "hierarchical sketch induction for paraphrase generation". In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020. Generating fact checking explanations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7352–7364, Online. Association for Computational Linguistics.

Tal August, Katharina Reinecke, and Noah A. Smith. 2022. Generating scientific definitions with controllable complexity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8298–8317, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz. 2022. A metrological perspective on reproducibility in NLP. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Shubham Agarwal, Anastasia Shimorina, and Ehud Reiter. 2021a. A systematic review of reproducibility research in natural language processing. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 381–393, Online. Association for Computational Linguistics.

Anya Belz, Simon Mille, and David Howcroft. 2020. Disentangling the properties of human evaluation methods: A classification system to support comparability, meta-evaluation and reproducibility testing. In *Proceedings of the 13th International Conference on Natural Language Generation*.

Anya Belz, Anastasia Shimorina, Shubham Agarwal, and Ehud Reiter. 2021b. The reprogen shared task on reproducibility of human evaluations in nlg: Overview and results. *INLG 2021*, page 249.

Anya Belz, Anastasia Shimorina, Maja Popovic, and Ehud Reiter. 2022. The 2022 reprogen shared task on reproducibility of evaluations in nlg: Overview and results. *INLG 2022*, page 43.

Anya Belz and Craig Thomson. 2023. The 2023 ReproNLP shared task on reproducibility of evaluations in NLP: Overview and results. In *Proceedings of the 3rd Workshop on Human Eval-*

*uation of NLP Systems*, pages 35–48, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *Proceedings of the Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

António Branco, Nicoletta Calzolari, Piek Vossen, Gertjan Van Noord, Dieter van Uytvanck, João Silva, Luís Gomes, André Moreira, and Willem Elbers. 2020. A shared task of a new, collaborative type to foster reproducibility: A first exercise in the area of language science and technology with REPROLANG2020. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 5539–5545, Marseille, France. European Language Resources Association.

Thiago Castro Ferreira, Diego Moussallem, Ákos Kádár, Sander Wubben, and Emiel Krahmer. 2018. NeuralREG: An end-to-end approach to referring expression generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1959–1969, Melbourne, Australia. Association for Computational Linguistics.

Bharathi Raja Chakravarthi, Navya Jose, Shardul Suryawanshi, Elizabeth Sherly, and John Philip McCrae. 2020. A sentiment analysis dataset for code-mixed Malayalam-English. In *Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL)*, pages 177–184, Marseille, France. European Language Resources association.

Tanvi Dinkar, Gavin Abercrombie, and Verena Rieser. 2024. Dexpert evaluation? reproducing human judgements of the fluency of generated text. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Xiachong Feng, Xiaocheng Feng, Libo Qin, Bing Qin, and Ting Liu. 2021. Language model as an annotator: Exploring DialoGPT for dialogue summarization. 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 1479–1491, Online. Association for Computational Linguistics.

Andra-Maria Florescu, Marius Micluta-Campeanu, and Liviu P. Dinu. 2024. Once upon a replication: It is humans' turn to evaluate ai's understanding of children's stories for qa generation. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Vivian Fresen, Mei-Shin Wu-Urbanek, and Steffen Eger. 2024. Humeval 24 reproduction report for paper 0043: Language model as an annotator: Exploring dialogpt for dialogue summarization. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Saadia Gabriel, Skyler Hallinan, Maarten Sap, Pemi Nguyen, Franziska Roesner, Eunsol Choi, and Yejin Choi. 2022. Misinfo reaction frames: Reasoning about readers' reactions to news headlines. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3108–3127, Dublin, Ireland. Association for Computational Linguistics.

Mingqi Gao, Jie Ruan, and Xiaojun Wan. 2024. A reproduction study of the human evaluation of the coverage of fact checking explanations. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Sebastian Gehrmann, Elizabeth Clark, and Thibault Sellam. 2023. Repairing the cracked foundation: A survey of obstacles in evaluation practices for generated text. *J. Artif. Int. Res.*, 77.

Javier González Corbelle, Ainhoa Vivel Couso, Jose Maria Alonso-Moral, and Alberto Bugarín-Diz. 2024. Reproducing the human evaluation of the dexperts controlled text generation method. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Yuxuan Gu, Xiaocheng Feng, Sicheng Ma, Lingyuan Zhang, Heng Gong, and Bing Qin. 2022. A distributional lens for multi-aspect controllable text generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1023–1043, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Yuxuan Gu, Xiaocheng Feng, Sicheng Ma, Lingyuan Zhang, Heng Gong, Weihong Zhong, and Bing Qin. 2023. Controllable text generation via probability density estimation in the latent space. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12590–12616, Toronto, Canada. Association for Computational Linguistics.

Tom Hosking, Hao Tang, and Mirella Lapata. 2022. Hierarchical sketch induction for paraphrase generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2489–2501, Dublin, Ireland. Association for Computational Linguistics.

David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 169–182, Dublin, Ireland. Association for Computational Linguistics.

Zdeněk Kasner and Ondrej Dusek. 2022. Neural pipeline for zero-shot data-to-text generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3914–3932, Dublin, Ireland. Association for Computational Linguistics.

Filip Klubička and John Kelleher. 2024. Reprohum #1018-09: Reproducing human evaluations of redundancy errors in data-to-text systems. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Mateusz Lango, Patricia Schmidtova, Simone Balloccu, and Ondrej Dusek. 2024. Evaluating summarization models: investigating the impact of education and language proficiency on reproducibility. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Yiru Li, Huiyuan Lai, Antonio Toral, and Malvina Nissim. 2024. Report for reprohum project 0033. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. DExperts: Decoding-time controlled text generation with experts and anti-experts. 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 6691–6706, Online. Association for Computational Linguistics.

Tyler Loakman and Chenghua Lin. 2024. Human evaluation reproduction report for generating fact checking explanations. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Michela Lorandi and Anya Belz. 2024. Reproducing the metric-based evaluation of a set of controllable text generation techniques. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Saad Mahamood. 2024. Reproducing human evaluations of end-to-end approaches to referring expression generation. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Zola Mahlaza, Toky Raboanary, Kyle Seakgwa, and C. Maria Keet. 2024. Another evaluation of readers' reactions to news headlines. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Irene Mondella, Huiyuan Lai, and Malvina Nissim. 2024. Report for reprohum project 0892-01. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Sachin Sasidharan Nair, Tanvi Dinkar, and Gavin Abercrombie. 2024. Exploring reproducibility of human-labelled data for code-mixed sentiment analysis. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Matthew Shardlow and Raheel Nawaz. 2019. Neural text simplification of clinical letters with a domain specific phrase table. In *Proceedings of the*

*57th Annual Meeting of the Association for Computational Linguistics*, pages 380–389, Florence, Italy. Association for Computational Linguistics.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Uri Simonsohn. 2015. Small telescopes: Detectability and the evaluation of replication results. *Psychological Science*, 26(5):559–569. PMID: 25800521.

Chris van der Lee, Albert Gatt, Emiel van Miltenburg, and Emiel Krahmer. 2021. Human evaluation of automatically generated text: Current trends and best practice guidelines. *Computer Speech & Language*, 67:101151.

Emiel van Miltenburg, Anouck Braggaar, Nadine Braun, Martijn Goudbeek, Emiel Krahmer, Chris van der Lee, Steffen Pauws, and Frédéric Tomas. 2024. How reproducible are fluency ratings of generated text? a reproduction of august et al. 2022. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Lewis Watson and Dimitra Gkatzia. 2024. Reprohum #0712-01: Reproducing human evaluation of meaning preservation in paraphrase generation. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*, Turin, Italy. Association for Computational Linguistics.

Bingsheng Yao, Dakuo Wang, Tongshuang Wu, Zheng Zhang, Toby Li, Mo Yu, and Ying Xu. 2022. It is AI's turn to ask humans a question: Question-answer pair generation for children's story books. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 731–744, Dublin, Ireland. Association for Computational Linguistics.

## A. The ReproHum Common Approach to Reproduction

In order to ensure comparability between studies, we agreed the following common-ground approach to carrying out reproduction studies:

1. Plan for repeating the original experiment in a form that is as far as possible identical to the original experiment, ensuring you have all required resources in place, then apply to research ethics committee for approval. If any aspect of the original experiment is unclear, contact the ReproHum coordinator who will either obtain clarification from the author, or create a sensible design that will then be used by all partner labs reproducing that experiment.

2. If participants were paid during the original experiment, determine pay in accordance with the ReproHum common procedure for calculating fair pay (Belz et al., 2023).

3. Following ethical approval start the reproduction study following the steps below. Contact the ReproHum team with any questions rather than the original authors, as they have already provided us with all the resources and information they have. Don't communicate with other ReproHum teams about their reproduction studies. This is to avoid inadvertently affecting outcomes.

4. Complete HEDS datasheet.

5. Identify the following types of results reported in the original paper for the experiment:

   (a) Type I results: single numerical scores, e.g. mean quality rating, error count, etc.

   (b) Type II results: sets of numerical scores, e.g. set of Type I results .

   (c) Type III results: categorical labels attached to text spans of any length.

   (d) Qualitative conclusions/findings stated explicitly in the original paper.[6]

6. Carry out the allocated experiment exactly as described in the HEDS sheet.

7. Report the results in the following form:

   (a) Description of the original experiment.

   (b) Description of any differences in your repeat experiment.

   (c) Side-by-side presentation of all results (8a-d above) from original and repeat experiments, in tables.

   (d) Report quantified reproducibility assessments as follows:

      i. Type I results: Small-sample oefficient of variation CV* (Belz, 2022).

      ii. Type II results: Pearson's r, Spearman's $\rho$.

      iii. Type III results: Multi-rater: Fleiss's $\kappa$; Multi-rater, multi-label: Krippendorff's $\alpha$.

---

[6]We now call these Type IV results.

iv. Conclusions/findings: Side-by-side summary of conclusions/findings that are / are not confirmed in the repeat experiment.

## B. Rounding Policy

The python script used to calculate results uses HALF_UP rounding rather than the python default of bankers rounding. Numbers are only ever rounded at the stage of presentation, i.e., the full-precision CV* values are used to calculated the means, rather than the 2 decimal place ones.

For Pearson and Spearman correlations we never round up from 0.99 in order to avoid giving the impression of a perfect correlation where one does not exist.

# Once Upon a Replication: It is Humans' Turn to Evaluate AI's Understanding of Children's Stories for QA Generation

**Andra-Maria Florescu[1,*], Marius Micluța-Câmpeanu[1,*], Liviu P. Dinu[2]**

[1]Interdisciplinary School of Doctoral Studies
[2]Faculty of Mathematics and Computer Science
University of Bucharest, Romania
{andra-maria.florescu,marius.micluta-campeanu}@s.unibuc.ro
ldinu@fmi.unibuc.ro

### Abstract

The following paper presents the outcomes of a collaborative experiment on human evaluation from the ReproNLP 2024 shared task, track B, part of the ReproHum project. For this paper, we evaluated a QAG (question-answer generation) system centered on English children's storybooks that was presented in a previous research, by using human evaluators for the study. The system generated relevant QA (Question-Answer) pairs based on a dataset with storybooks for early education (kindergarten up to middle school) called FairytaleQA. In the framework of the ReproHum project, we first outline the previous paper and the reproduction strategy that has been decided upon. The complete setup of the first human evaluation is then described, along with the modifications required to replicate it. We also add other relevant related works on this subject. In conclusion, we juxtapose the replication outcomes with those documented in the cited publication. Additionally, we explore the general features of this endeavor as well as its shortcomings.

**Keywords:** ReproNLP, QAG system, FairytaleQA, Reproduction

## 1. Introduction

In the field of Natural Language Processing (NLP), reproducibility is crucial for democratizing and understanding better the mechanisms of the field (Storks et al., 2023). Nevertheless, there are still issues and no widely recognized, appropriate procedure for carrying out replications of earlier research. A major factor that continues to make reproduction challenging to accomplish is the evaluations conducted by both humans and computers (Belz et al., 2023a; Pineau et al., 2021). A wide range of variables, such as imprecise data, incorrect experiments, and disagreement among the human assessors, make human evaluation one of the major obstacles to accurately replicating previous research (Thomson et al., 2024; Belz et al., 2023b; Popović, 2021).

The present study focuses on human evaluation of prior NLP research and is part of the ReproNLP 2024 shared task on Reproducibility of Evaluations in NLP (Belz and Thomson, 2024), namely on the Track B task associated with the ReproHum project. The plan was to undertake the study again and try to replicate the findings. For this project, we replicated an NLP study in which we evaluated a QAG (Question and Answer Generation) system conducted by Yao et al. (2022) and compared the outcomes of this replication to the original findings. To our knowledge, the present study represents the first attempt of replicating these results.

Section 2 focuses on presenting the original study, QAG systems, the common strategy for evaluating QAG systems, and related studies presented in section 3. Section 4 explains how the NLP evaluation was replicated. It begins by outlining the contents of the selected paper and then goes into depth about every aspect of the evaluation that was replicated. Section 5 presents and discusses the findings from the replicated evaluation concerning the original study. Lastly, Section 7 offers some closing thoughts and future works related to this project.

In alignment with open science principles, we make available all code and data employed in this investigation for the benefit of the scientific community and future research endeavors[1].

## 2. QAG system

The original study, "It is AI's Turn to Ask Humans a Question: Question-Answer Pair Generation for Children's Story Books" (Yao et al., 2022) examined the question-answer pair generation task (QAG) in the context of early childhood education (kindergarten through middle school). The original study implemented a QA-pair generation pipeline, which, as observed in human and automated evaluation, effectively supported the objective of automatically generating high-quality questions and

---

[*]These authors contributed equally to this work.

---

[1]https://github.com/mcmarius/ReproNLP-2024

answers at scale. This was achieved by leveraging a newly-constructed expert-annotated QA dataset built upon child-oriented fairy tale storybooks (FairytaleQA, Xu et al., 2022).

Five non-native English speakers were selected for the study's human evaluation in order to assess the QAG system's capacity to produce high-quality Question-Answer pairs. Furthermore, in addition to the QAG system, the ground truth and the PAQ system (Lewis et al., 2021) were evaluated by human evaluators who were blind to the system they were assessing. Ground truth QA pairs were written by human annotators of the FairytaleQA dataset, while the PAQ system consists of two components: a passage selection model and an answer extraction model. The PAQ system is supported by the PAQ dataset, a corpus of 65 million QA pairs that were automatically generated.

Each QA pair's model of origin was unknown to the participants. Using a five-point Likert scale, the participant was asked to rate the QA pairs along three dimensions:

- Readability: The generated QA pair is in readable English grammar and words.

- Question Relevancy: The generated question is relevant to the storybook section.

- Answer Relevancy: The generated answer is relevant to the question.

According to the original paper, seven novels were chosen at random, and then ten sections from those seven books were chosen randomly (for a total of seventy-QA pairings). To ensure coding consistency, each participant was asked to rate these identical 70 QA pairs. After this step, ten books (five from the test and five from the validation divides) were then chosen at random, and four sections from each book were chosen randomly once again. There are, on average, nine QA-pairs in each section (three for each model). Two coders were assigned at random to each section. Overall, each coder coded four volumes, or sixteen sections and about 140 QA-pairs. A total of 722 QA-pairs were scored. T-tests were also used in the original study to determine if the difference between models is statistically significant.

## 3. Related works

The goal of automatic question generation, or QG, is to extract meaningful questions and desired responses from text sections. In the past, rule-based or neural models were employed; new developments have made neural models more popular. These models—sequence-to-sequence models in particular—are capable of creating excellent questions by utilizing prior knowledge and anticipated responses. However, their value is limited because they frequently require another system in order to obtain the correct answers. Additionally, there aren't many publicly available data sets for QG systems that may generate both questions and answers. An alternative method concentrates on teaching QG models solely on context, enabling them to produce distinct question kinds for varying text lengths. State-of-the-art (SOTA) systems use pre-trained language models (PLMs) like Google T5 and GPT-3 for instructional neural question generation[2]. These pre-trained models on large-scale text corpora allow for the creation of questions with zero effort and no further training. GPT models have the ability to generate educational questions, as a recent study has shown (Bulathwela et al., 2023). Therefore, since the first study was published, numerous additional studies have been carried out on QA AI systems, some of which have been especially focused on issues related to education. Ushio et al. (2023) released `AutoQG`, a multilingual web-based quality assurance system, and `lmqg`, a Python module for QA generation, fine-tuning, and evaluation. This user-friendly code might be advantageous to both developers and end users who require customized models or fine-grained controls for development.

## 4. Reproduction of the human evaluation

We were given a document by the task organizers with more details regarding the human evaluation procedure to help with our reproduction experiment, even though we weren't able to interact with the authors directly. The ReproNLP 2024 project team corresponded with the authors to obtain this information prior to initiating the ReproNLP 2024 shared task. The document covers details regarding the task configuration given to the human evaluators, including the methods used. In addition, we fulfilled the experiment's requirements by filling out a Human Evaluation Datasheet (HEDS, Shimorina and Belz, 2022)[3]. This form consists of details on the assignment, the evaluators' characteristics, and the gathered annotation information from them.

For this paper, we attempted to follow the original procedures given by the prior study and the extra information obtained as closely as possible in order to replicate the human evaluation. Five human subjects were employed in the initial study to rate each of the three QA systems. The system outputs were selected randomly by the authors of the

---

[2]This was the state-of-the-art when the initial study was published in 2022.

[3]The HEDS document is available here: https://github.com/nlp-heds/repronlp2024

original paper, so we used the same set of examples. To ensure our reproducibility, we aimed for the same number of evaluators. We first posted our requirements in an announcement sent to the student representative who shared it on their communication channels. The only inclusion criterion was for the student to be at least in 3rd year. As a result, five undergraduate male BSc and BEng students who are not native English speakers, but speak it fluently answered our request. Gender was not a criterion used for selection, other students could have participated as well. They received no monetary compensation for their involvement, which was instead taken into account as part of their educational curriculum practice hours for which they needed academic credits.

To enable our evaluators to score the narrative sections and QA pairings from the three systems—Ground truth, PAQ, and the original paper's system (called "Ours")—blindly on a scale of 1 to 5 for readability and relevance for questions and answers, the students were each given an Excel file with 7 columns: id (internal), section text, question text, answer text and 3 columns corresponding to each of the 3 ratings they have to provide. The students annotated at their own pace from their place of choosing, but they were instructed that they had a deadline of one week to complete the Excel sheets with their evaluations. Each student read the sections, questions and answers not knowing what QA system they were assessing as well as not being aware of the other annotators in order to have an unbiased evaluation. They reported that on average their annotation took up to 5 hours. We had no pre-coding training or detailed coding guidelines as indicated in the document received from the ReproNLP task organizers regarding the original study.

## 5. Findings

In this section, we present the main outcomes of our study, along with qualitative and quantitative analyses that strive to explain the disparities from the preceding research. We show our approach for determining the inter-annotator agreement, along with the differences in statistical significance of the results between the two experiments. We also include a quantified reproducibility assessment (QRA) (Belz et al., 2022).

### 5.1. Inter-annotator agreement

First, we attempt to compute the inter-coder reliability score (Krippendorff's alpha, Krippendorff, 2011) for both experiments. If we assume that the authors limited the scope of their pre-coding stage to ground truth examples, we are able to partially con-

firm the claim from the original paper that shows a high level of agreement between all annotators. We determined this agreement based on the available data that we received.

Given that each sample is coded only by two raters, we compute the overall agreement by averaging the individual agreements between each pair of raters with common examples[4]. There are 3 systems to be evaluated, 5 annotator pairs and 3 evaluation dimensions, leading to a total of 45 individual pairs. Out of these 45 pairs, there are 12 instances with acceptable alpha values over 0.67, resulting in an overall Krippendorff's alpha score of 0.43 for the initial experiment. We note that 9 out of those 12 instances are for ground truth examples. Only 3 out of 30 pairs show an alpha value over 0.67 for their system and the PAQ system, all of them for answer relevancy. A breakdown of these values by system and evaluation dimension is shown in Table 1.

We rely on the only Krippendorff's alpha Python package that provides support for ordinal levels of measurement (Castro, 2017), since we need to distinguish between low and high score differences. Upon some investigations, we find that this implementation does not take into account situations with perfect or almost perfect agreement and only one conforming value (e.g. most ratings have a common score of 5, but there is no example where both labelers give a score of 4), leading to spurious values that erroneously entail no agreement when analyzing subsets of the original data. These issues could be (partially) mitigated by determining the inter-coder reliability score on a larger sample where other identical values are more likely to appear. If this is not feasible, researchers should at least properly specify the software packages used[5].

Next, we calculate Krippendorff's alpha for our evaluators, acknowledging that no pre-coding practice took place due to missing coding guidelines. From the total of 45 pairs, there are 8 instances with an alpha value over 0.67. The automated systems obtain alpha values over 0.67 in 5 out of the same 30 instances, with two additional alpha scores over 0.65. These results show a slightly higher agreement for the automated systems compared to the original paper. Again, this agreement is observed mostly for answer relevancy, with two instances being for question relevancy. The overall

---

[4]If we compute the overall agreement directly, the score underestimates the real agreement due to the sparsity of data.

[5]Additionally, we computed the agreement scores using the R package `irr`. While the case of identical values is handled correctly by `irr`, it does not allow specifying the domain of possible values. The results are identical, but we note that both implementations show no agreement in other corner cases of almost perfect agreement.

|  | Ours | PAQ Baseline | Groundtruth |
|---|---|---|---|
| Readability | 0.25 | 0.24 | 0.94 |
| Question relevancy | 0.38 | 0.33 | 0.35 |
| Answer relevancy | 0.45 | 0.44 | 0.45 |
| **Overall agreement** | | | **0.43** |

Table 1: Inter-annotator agreement measured by Krippendorff's alpha for the original paper for each system and evaluation dimension. Each cell shows an average of 5 pairs of coders.

|  | Ours | PAQ Baseline | Groundtruth |
|---|---|---|---|
| Readability | -0.06 | 0.05 | -0.13 |
| Question relevancy | 0.35 | 0.46 | 0.42 |
| Answer relevancy | 0.51 | 0.50 | 0.46 |
| **Overall agreement** | | | **0.27** |

Table 2: Inter-annotator agreement measured by Krippendorff's alpha for the replication experiment for each system and evaluation dimension. Each cell shows an average of 5 pairs of coders.

Krippendorff's alpha score is 0.27 for the replication study. Table 2 offers a systematic overview of the agreement by system and evaluation dimension, revealing marginally better agreements for relevancy scores than the initial paper.

## 5.2. Statistical significance of the results

After receiving the annotated files, we perform a sanity check for each evaluator by counting the number of samples for which ratings have any absolute difference in contrast with the original labels. This step reveals that one of our annotators assigned substantially inferior grades, prompting us to omit these biased scores from the statistical tests.

For completeness, the initial results are displayed in Table 3, while the same results ignoring the biased labeler are shown in Table 4. We first validate the assumptions of *t-tests* through Shapiro-Wilk tests, confirming that the scores for automated systems ("Ours" and PAQ) are normally distributed. As expected, the ground truth distribution is skewed since most ratings are 4 or above.

The proposed model ("Ours", avg = 4.52, s.d. = 0.79) significantly outperforms PAQ for the *Readability* dimension (avg = 4.13, s.d. = 1.04, $t(382)$ = 4.07, $p < 0.01$), albeit not as satisfactory as the ground-truth (avg = 4.67, s.d. = 0.55, $t(392)$ = -2.64, $p < 0.01$).

In terms of the *Question relevancy* dimension, ground-truth (avg = 4.77, s.d. = 0.71) surpasses the proposed model (avg = 3.92, s.d. = 1.37), which in turn is significantly better than the PAQ baseline (avg = 3.39, s.d. 1.60, $t(382)$ = 2.05, $p < 0.05$)[6].

Finally, for the *Answer relevancy* dimension, the ground-truth obtains by far the best ratings (avg = 4.58, s.d. = 0.92). Unlike the original paper, our experiments do not display a considerable distinction between the proposed model (avg = 3.39, s.d. = 1.60) and the PAQ model (avg = 3.42, s.d. 1.62, $t(382)$ = -0.22, $p$ = .82), though we confirm that this observation is not statistically significant.

We include the human evaluation results of the original paper in Table 5 to aid comparisons with our replication experiments, although we do not repeat the *t-tests* results from the initial paper here since we are able to confirm both the exact numbers and the results of the statistical tests.

## 5.3. Quantified Reproducibility Assessment

Quantified reproducibility assessment (QRA), introduced by Belz et al. (2022), aims to provide an impartial framework for determining the extent of reproducibility across different tasks and types of evaluation. This is achieved by computing a single score known as precision for each value of interest, which enables comparability between studies.

In accordance with the guidelines provided by the task organizers, we use the unbiased coefficient of variation (CV*) for small sample sizes (Belz, 2022) as a measure for precision. This score is determined independently for each of the three dimensions, as shown in Tables 6, 7 and 8, along with Pearson's correlation coefficient $r$ and Spearman's correlation coefficient $\rho$.

As mentioned in the previous section, the results

---

[6]If we include the biased labeler, the difference between "Ours" and PAQ is no longer significant: $t(478)$ = 1.79, $p$ = 0.07. This would be the only place where the biased labeler meaningfully affects the statistical tests.

|  | Ours | | PAQ Baseline | | Groundtruth | |
|---|---|---|---|---|---|---|
|  | M | SD | M | SD | M | SD |
| Readability | 4.52 | 0.75 | 4.17 | 1.22 | 4.71 | 0.52 |
| Question Relevancy | 3.83 | 1.30 | 3.61 | 1.35 | 4.71 | 0.73 |
| Answer Relevancy | 3.20 | 1.56 | 3.20 | 1.57 | 4.46 | 1.03 |

Table 3: Human evaluation results of the reproduction study

|  | Ours | | PAQ Baseline | | Groundtruth | |
|---|---|---|---|---|---|---|
|  | M | SD | M | SD | M | SD |
| Readability | 4.52 | 0.79 | 4.13 | 1.04 | 4.67 | 0.55 |
| Question Relevancy | 3.92 | 1.37 | 3.62 | 1.45 | 4.77 | 0.71 |
| Answer Relevancy | 3.39 | 1.60 | 3.42 | 1.62 | 4.58 | 0.92 |

Table 4: Human evaluation results of the reproduction study excluding the biased labeler

|  | Ours | | PAQ Baseline | | Groundtruth | |
|---|---|---|---|---|---|---|
|  | M | SD | M | SD | M | SD |
| Readability | 4.71 | 0.70 | 4.08 | 1.13 | 4.95 | 0.28 |
| Question Relevancy | 4.39 | 1.15 | 4.18 | 1.22 | 4.92 | 0.33 |
| Answer Relevancy | 3.99 | 1.51 | 3.90 | 1.62 | 4.83 | 0.57 |

Table 5: Human evaluation results of the original paper

| System | Orig | Repl | CV$^*$ | $r$ | $\rho$ |
|---|---|---|---|---|---|
| Ours | 4.71 | 4.52 | 4.10 | | |
| PAQ | 4.08 | 4.17 | 2.18 | 0.99 | 1 |
| GT | 4.95 | 4.71 | 4.95 | | |

Table 6: Precision metrics for the readability dimension showing the degree of reproducibility. CV$^*$ is computed using $n = 2$. Pearson's correlation and Spearman's correlation are denoted by $r$ and $\rho$ respectively. *Orig* indicates results from the initial experiment by Yao et al. (2022). *Repl* refers to replicated scores. *GT* represents ground truth scores.

| System | Orig | Repl | CV$^*$ | $r$ | $\rho$ |
|---|---|---|---|---|---|
| Ours | 4.39 | 3.83 | 13.58 | | |
| PAQ | 4.18 | 3.61 | 14.59 | 0.99 | 1 |
| GT | 4.92 | 4.71 | 4.35 | | |

Table 7: Precision metrics for the question relevancy dimension showing the degree of reproducibility. We use the conventions from Table 6.

are not statistically significant if we include the problematic labeler. We obtain $r = 0.99, p = 0.056$ and $\rho = 1, p = 0.0$ for *Readability*, $r = 0.99, p = 0.056$ and $\rho = 1, p = 0.0$ for *Question relevancy*, and

| System | Orig | Repl | CV$^*$ | $r$ | $\rho$ |
|---|---|---|---|---|---|
| Ours | 3.99 | 3.20 | 21.90 | | |
| PAQ | 3.90 | 3.20 | 19.66 | 0.99 | 0.87 |
| GT | 4.83 | 4.46 | 7.94 | | |

Table 8: Precision metrics for the answer relevancy dimension showing the degree of reproducibility. We use the conventions from Table 6.

$r = 0.99, p = 0.03$ with $\rho = 0.86, p = 0.33$ for *Answer relevancy*.

QRA results display low CV$^*$ values for readability, while relevancy scores showcase a substantial gap between QAG systems and ground truth, prompting the need for precise coding instructions.

## 5.4. Reproduction results

In order to talk about the differences between the original study and ours, we had online meetings with the five human evaluators, focusing on examples with conflicting scores when compared to the original labels. Together, we examined the Excel documents that they had annotated, and we asked them to justify the scores they had given for readability, question relevancy, and response relevancy. It appears that the majority of our annotators based their remarks on their personal interpretations of the

texts, language proficiency, comprehension, and instances where errors resulting from a failure to pay attention to the texts affected the scoring.

We conduct a quantitative analysis stemming from the findings recorded as part of the discussions with our annotators. We synthesize our interpretations for labeling discrepancies in Table 9, noting that we consider examples as belonging exclusively to one error category to better observe systematic mistakes.

One persistent problem with the PAQ system was that it would repeatedly replace the named entities in the questions with "val", for instance: "What did val give to the dead man?". Out of 240 samples (120 unique questions), "val" appears in 106 of them. This caused our annotators to assign readability scores that were lower than those in the original study for 19 occurrences.

The responses were incomplete in 11 cases, like the following:

> Question: What did the man give his son?
>
> Answer: falcon.

The complete answer here would have been "gun, dog and falcon."

Since we made a methodological error by only providing generic scoring instructions without specific restrictions or details, one labeler relied on simple heuristics and primarily assigned low values for single-word responses even if they were otherwise relevant and readable. These account for 19 QA mislabeled pairs. Still, we argue that an educational QA system should seek to include connectives and proper punctuation marks as part of their answers. For example:

> Question: What weapon did val use to cut down a tree?
>
> Answer: axe

Similarly, high ratings were provided for questions or answers that resemble verbatim portions of the story, despite the lack of meaning or importance. The previous article's QA system ("Ours") tends to generate such copy-paste fragments from the story sections, in some cases being illegible: "the son - in - law ate nothing though his wife 's parents , with kind words and friendly gestures , kept urging him to help himself".

It should be noted that the initial labeling is also prone to human errors. These situations are infrequent, but they represent more than 10% of divergent ratings. The following QA pair has received marks of 4 and 5 for readability in the previous experiment, despite the nonsensical nature:

> Question: What kind of garlic would a cow be good for?
>
> Answer: garlic.

Our annotators disregarded the possibility that some questions and answers were pertinent and might have been inferred from the sections, thus focusing only on explicit textual matches. We also noticed that regarding readability, the scoring was influenced by the QA pair, although there were instances in which the question was readable, while the answer was not, thus influencing the rating. We suggest that readability should also be scored independently for question and answer.

## 6. Discussion

We first reiterate the contributions of the original work and the extent of our replication before discussing the consequences of our findings. In order to enhance the accuracy of automated question-answer generation systems in educational contexts—specifically children's storybooks—Yao et al. (2022) introduced an innovative technique. They demonstrated the superiority of their approach over current state-of-the-art models on two datasets, PAQ and 2-step baseline systems, as well as ground truth (human educational experts), using a combination of automatic and human evaluation approaches. Our replication was restricted to the human evaluation task described in their study, which assessed the produced questions and answers for readability, relevance of the questions, and relevance of the responses in relation to the story's segments.

As stated by Arvan and Parde (2023) in their reproducibility article from ReproNLP 2023, there was insufficient information in the research paper to replicate the original human evaluation in its entirety. This is likely due to the fact that, in the current research climate, NLP research is too focused on novelty and format compliance, rather than providing a clear explanation of the methodologies used.

Given that human evaluation is carried out by humans, personality, culture, expertise, and comprehension can all lead to significant biases (Amidei et al., 2018). This is why, in order to minimize errors made by humans as much as possible, explicit standards for the evaluations are required to obtain less ambiguous interpretations of the annotators. For example, regarding this study and others that focus on QA systems, evaluation dimensions such as readability should be assessed separately for questions and answers.

## 7. Conclusions

All in all, we managed to replicate the original study. However, our annotators considered answer length, which affected their low scoring because of a methodological error on our part. Furthermore, some questions might have been legitimate even

| Error category | Count |
|---|---|
| Readability | 14 |
| Incomplete question | 1 |
| Irrelevant question | 14 |
| Incomplete answer | 11 |
| Right answer in another context | 2 |
| Wrong answer | 10 |
| Short answer | 19 |
| Perception, comprehension | 14 |
| "val" mentioned | 19 |
| Methodological errors | 6 |
| Human error (reproduction study) | 13 |
| Human error (original study) | 16 |
| **Total** | **139** |

Table 9: Quantitative analysis of divergent answers with an absolute score difference of 3 or 4 in at least one dimension

though the answer was only inferred rather than explicitly stated in the text; nonetheless, our labelers focused solely on information that was explicitly mentioned in the text, which led to another lower score than the original study.

As mentioned in several studies centered on human evaluation (Amidei et al., 2018), one's personality, language knowledge as well as own writing style influence drastically the scoring. This was present in our study as well. After discussing with our annotators, we noticed that in most cases, their personality and English understanding knowledge influenced the scoring. It is clear from comparing the two studies that there are some differences between the replicated and original results. The differences between the first study and ours are more likely to be the result of methodological errors because we were not given access to the entire set of original guidelines, as well as human errors made by the annotators in terms of comprehending and interpreting the assignments.

## 8. Limitations

Unlike the original research, we only employed BSc and BEng students for this study, and they came from a different field than the original work for the human evaluation. We took the most of the scant information available because we lacked the precise guidelines from the previous research hence having some methodological errors for the evaluations.

## 10. Bibliographical References

Jacopo Amidei, Paul Piwek, and Alistair Willis. 2018. Rethinking the agreement in human evaluation tasks. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 3318–3329, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Mohammad Arvan and Natalie Parde. 2023. Human evaluation reproduction report for data-to-text generation with macro planning. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 89–96, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Anya Belz. 2022. A metrological perspective on reproducibility in nlp. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Maja Popovic, and Simon Mille. 2022. Quantified reproducibility assessment of NLP results. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16–28, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz, Maja Popović, Ehud Reiter, Craig Thomson, and João Sedoc, editors. 2023a. *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*. INCOMA Ltd., Shoumen, Bulgaria, Varna, Bulgaria.

Anya Belz and Craig Thomson. 2024. The 2024 ReproNLP shared task on reproducibility of evaluations in NLP: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes,

Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023b. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *The Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Sahan Bulathwela, Hamze Muse, and Emine Yilmaz. 2023. Scalable educational question generation with pre-trained language models. In *Artificial Intelligence in Education*, pages 327–339, Cham. Springer Nature Switzerland.

Santiago Castro. 2017. Fast Krippendorff: Fast computation of Krippendorff's alpha agreement measure. https://github.com/pln-fing-udelar/fast-krippendorff.

Klaus Krippendorff. 2011. Computing Krippendorff's alpha-reliability.

Joelle Pineau, Philippe Vincent-Lamarre, Koustuv Sinha, Vincent Lariviere, Alina Beygelzimer, Florence d'Alche Buc, Emily Fox, and Hugo Larochelle. 2021. Improving reproducibility in machine learning research (a report from the neurips 2019 reproducibility program). *Journal of Machine Learning Research*, 22(164):1–20.

Maja Popović. 2021. Agree to disagree: Analysis of inter-annotator disagreements in human evaluation of machine translation output. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 234–243, Online. Association for Computational Linguistics.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Shane Storks, Keunwoo Yu, Ziqiao Ma, and Joyce Chai. 2023. NLP reproducibility for all: Understanding experiences of beginners. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10199–10219, Toronto, Canada. Association for Computational Linguistics.

Craig Thomson, Ehud Reiter, and Anya Belz. 2024. Common Flaws in Running Human Evaluation Experiments in NLP. *Computational Linguistics*, pages 1–11.

Asahi Ushio, Fernando Alva-Manchego, and Jose Camacho-Collados. 2023. A practical toolkit for multilingual question and answer generation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 86–94, Toronto, Canada. Association for Computational Linguistics.

Bingsheng Yao, Dakuo Wang, Tongshuang Wu, Zheng Zhang, Toby Li, Mo Yu, and Ying Xu. 2022. It is AI's turn to ask humans a question: Question-answer pair generation for children's story books. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 731–744, Dublin, Ireland. Association for Computational Linguistics.

## 11. Language Resource References

Lewis, Patrick and Wu, Yuxiang and Liu, Linqing and Minervini, Pasquale and Küttler, Heinrich and Piktus, Aleksandra and Stenetorp, Pontus and Riedel, Sebastian. 2021. *PAQ: 65 Million Probably-Asked Questions and What You Can Do With Them*. MIT Press.

Xu, Ying and Wang, Dakuo and Yu, Mo and Ritchie, Daniel and Yao, Bingsheng and Wu, Tongshuang and Zhang, Zheng and Li, Toby and Bradford, Nora and Sun, Branda and Hoang, Tran and Sang, Yisi and Hou, Yufang and Ma, Xiaojuan and Yang, Diyi and Peng, Nanyun and Yu, Zhou and Warschauer, Mark. 2022. *Fantastic Questions and Where to Find Them: FairytaleQA – An Authentic Dataset for Narrative Comprehension*. Association for Computational Linguistics.

# Exploring Reproducibility of Human-Labelled
# Data for Code-Mixed Sentiment Analysis

**Sachin Sasidharan Nair, Tanvi Dinkar, Gavin Abercrombie**
Heriot-Watt University, Edinburgh, Scotland
{ss2246, t.dinkar, g.abercrombie}@hw.ac.uk

## Abstract

Growing awareness of a 'Reproducibility Crisis' in natural language processing (NLP) has focused on human evaluations of generative systems. While labelling for supervised classification tasks makes up a large part of human input to systems, the reproduction of such efforts has thus far not been been explored. In this paper, we re-implement a human data collection study for sentiment analysis of code-mixed Malayalam movie reviews, as well as automated classification experiments. We find that missing and under-specified information makes reproduction challenging, and we observe potentially consequential differences between the original labels and those we collect. Classification results indicate that the reliability of the labels is important for stable performance.

**Keywords:** Reproducibility, Human Data Collection, Sentiment Analysis, Malayalam

## 1. Introduction

There has recently been growing awareness of a 'Reproducibility Crisis' in natural language processing (NLP) (Belz et al., 2021). This has focused on the apparent impossibility of reproducing human evaluation studies of the outputs of natural language generation (NLG) systems (Belz et al., 2023; Thomson et al., 2024). However, while text labelling makes up the largest part of human input to NLP projects, there have been almost no reported attempts (to our knowledge) to reproduce human label collection for NLP tasks outwith NLG evaluation.

In this study, part of ReproNLP[1] Track A (Belz and Thomson, 2024), we focus on one of the most active areas of NLP over the last two decades, sentiment analysis (Mäntylä et al., 2018): an NLP task that aims to categorise the sentiment expressed in textual data (Liu, 2012).

In addition, we delve into the complexities introduced to sentiment analysis by using code-mixed language data. We re-examine the Malayalam-English corpus of Chakravarthi et al. (2020), classified with one of five distinct labels (*Positive*, *Negative*, *Neutral*, *Mixed feelings*, and *Non-Malayalam*). We assess the challenges faced while re-annotating the original corpus, and while reproducing the processes followed by the original study including the annotation process and the re-implementation of automated classifiers, verifying whether we are able to achieve similar results to those of the original study.

An example item from the corpus is shown here:

> Ufff vere level ikkaaa ingha pwoli aahn
> *Another level, ikka you are awesome*
> Assigned Label: *Positive*

## 2. Background & Related Work

### 2.1. Reproducbility in NLP

In response to the widespread reproduciblity issues uncovered in other scientific fields (Baker, 2016), there have been increasing efforts to establish standards for reproducibility in NLP, such as workshops that aim to tackle these problems (e.g. and Machine Learning Reproducbility Challenge (MLRC), HumEval)(Belz et al., 2021). Other initiatives, such as reproducibility checklists[2] have been adopted at major conferences such as EMNLP and AAAI to foster the integrity and validity of experiments.

The conversation around reproducibility is nuanced. In their review, Belz et al. (2021) note that the definitions provided by six different sources had varied interpretations of reproducibility and replicability, lacking standardised definitions. This diversity further complicates the efforts to establish consistent reproduciblity practices in the field of NLP. Moreover, the discussions of Rougier et al. (2017) and Wieling et al. (2018) highlight the need for a common understanding that also involves transparency and openness to guide reproducibility efforts.

Gundersen and Kjensmo (2018) evaluated 400 research papers from major conferences IJCAI and AAAI revealing a lack of comprehensive documentation. Only an average of 20% to 30% of necessary variables were documented, which indicated a significant gap. Although there was a slight improvement in documentation over time, the reproducibility scores generally decreased as documentation requirements grew. This analysis confirms that lack of documentation is a significant

---

challenge faced in AI research reproducibility.

As this field progresses, it is clear that we must focus on addressing the challenges by the use of continued dialogue and action to develop standards for reproducibility. In this paper, we begin to shift the focus of NLP reproducibility research from NLG evaluations to data labelling for a supervised NLP classification task.

## 2.2. Understanding Sentiment in Code-Mixed Data

A major challenge for sentiment analysis is the ambiguities that hinder accurate classification. The classes and definitions can vary widely, which complicates the standardisation of this task across studies. Context and language can significantly affect the sentiment perceived. For example, Moore and Rayson (2018) show that identifying idioms, detecting sarcasm, and understanding the role of modifiers can influence the sentiment and the accuracy and replicability of the task. Typically, classifiers are designed only to process text that are written in high resource languages such as English. However, many other languages are used for digital communication (which is often the focus of a sentiment analysis task) and code-mixing is widely used in multilingual societies.

Malayalam (ml) is a Dravidian language[3] distinguished by its complex and rich phonetic and grammatical structures. Some research has been conducted on sentiment analysis for Malayalam text (Nair et al., 2014), for example focusing on tweets (Soumya and Pramod, 2020). Less research has been conducted on Malayalam-English code-mixed data, and a plausible reason for this is the lack of data availability. However, Chakravarthi et al. (2020) presented a corpus for sentiment analysis of code-mixed text for Malayalam-English, which we re-examine here.

The influence of code-mixing on annotator agreement and reproducibility has also received very little attention. One broadly related work is that of Abercrombie et al. (2023), who examined the impact of two factors, time and second language, on the inter- and intra-annotator agreement in German and English texts for a hate speech labelling task. Importantly, they found that label collection is not as repeatable as assumed even with the same annotators (in either language), which raises interesting questions on the reproducibility of multi-lingual data in general. In this study, we focus on what we believe is an understudied aspect of reproducibility in NLP, i.e. reproducibility of a sentiment analysis task using code-mixed data.

---

[3]For an overview of the complex linguistic landscape of South Asia, including Dravidian languages, see Hock and Bashir (2016).

## 3. Chakravarthi et al. (2020): A Sentiment Analysis Dataset for Code-Mixed Malayalam-English

The original study by Chakravarthi et al. (2020) is comprised of data collection and labelling, as well as automated classification experiments. We provide a brief overview here.

### 3.1. Original Data Collection

In the original study, Chakravarthi et al. (2020) extracted 116,711 sentences from comments posted on YouTube about trailers for Malayalam movies from the year 2019, using the search term '*Malayalam movie 2019*', excluding instances that were in Malayalam script. The data gathered was then filtered to exclude any data that was non-code-mixed, i.e purely in English. The code-mixed content was then preprocessed, specifying that emojis were removed from sentences and sentences exceeding more than 15 words or fewer than 5 words were discarded. The resulting corpus contains 6,738 instances.

### 3.2. Original Annotation Process

The initial data labelling process was carried out by volunteer annotators. The label schema consisted of the following labels: *positive*, *negative*, *mixed feelings*, *neutral* and *non-Malayalam*. We follow the annotation process detailed in the original study which consists of three steps:

- **First Step:** Each item was labelled by two annotators independently. Items with the same two labels were considered finalised.

- **Second Step:** Items with label disagreements were annotated by a third annotator. Where agreement could be found among the three annotators, the labels were decided by majority vote (i.e. two out of three labels).

- **Third Step:** If there was no majority, these items were subsequently reviewed by two other annotators. Labels were again decided by majority vote.

As well as the three steps mentioned above, the original study omits to mention what was done with sentences still having label disagreements following the third step. These samples could have all five labels differing, or an absence of a majority label (i.e., two votes to two labels respectively and one vote to another). On enquiry, the authors responded that they discarded items on which there was no agreement after all three stages, and these are not included in the data that is made available in the original study.

### 3.3. Original Classification Experiments

Chakravarthi et al. (2020) used a range of ML classifiers such as logistic regression (LR), support vector machines (SVM), Random Forests (RF), K-nearest neighbours (KNN), and multinominal naive Bayes (MNB) along with Term-Frequency Inverse Document Frequency (TF-IDF) for feature selection. Additionally, they also implemented four deep learning classifiers: 1D Dimensional Convolution (Zhou et al., 2016), Dynamic Meta-Embeddigs (DME), Contextualised DME (CDME) (Kiela et al., 2018) and BERT (Devlin et al., 2019). We select the two classifiers that attained highest performance (as detailed in Table 7), LR and BERT, for our classification reproduction study (see section 5).

## 4. Reproduction Study

### 4.1. Data

To maximise our resources, we selected only the test set split of the corpus for re-annotation.

### 4.2. Annotation Reproduction

We endeavoured to follow the original data annotation process as far as possible. However, to avoid discarding further items that were included in the original dataset, we incorporated a slight modification. For sentences that did not reach consensus among the initial five annotators (i.e., until the third step as described in subsection 3.2), rather than exclude such data, we added two further steps:

- **Fourth Step:** In scenarios where there was no clear majority, a sixth annotator was introduced to label the remaining items with label disagreements.

The reproduction study in theory should not elicit high disagreement amongst the annotators, as samples that did not have a majority label have already been discarded from the available data of the original study. Surprisingly, our re-annotation still yielded 61 items with label disagreements unresolved after step three, and 21 after step four. We therefore added a fifth step:

- **Fifth Step:** In case of unresolved disagreement among six annotators, such as an even split across three labels, ties between two labels, or a distributed disagreement (e.g., 2 *Positive*, 2 *Negative*, 1 *Neutral*, 1 *Mixed feelings*), the remaining 21 label disagreements were resolved by the first author of this work.

This process is detailed in Table 1.

| Step | No. of Annotators | Labels collected | Dis- agreements |
|---|---|---|---|
| 1 | 2 | 2 | 592 |
| 2 | 1 | 3 | 175 |
| 3 | 2 | 5 | 61 |
| 4 | 1 | 6 | 21 |
| 5 | 1 | 7 | 0 |

Table 1: Numbers of annotators, labels collected per item, and disagreements for each step.

### 4.3. Annotation Platform

We collected labels using the application MS Excel. The data was split into batches and was populated for annotation, which was undertaken on the annotators' personal computers.

### 4.4. Annotation Instructions

Chakravarthi et al. (2020) provide label definitions which they loosely adapt from Mohammad (2016), but little detail of the actual instructions given to the annotators. We contacted the first author for further clarification, but these were unavailable. We used the original study's label categories and definitions, and added an additional objective to instruct annotators on the task:

- **Objectives:** Categorise each sentence into one of the following segments.

    - **Positive:** There is an explicit or implicit clue in the text suggesting that the speaker is in a positive state, i.e., happy, admiring, relaxed, and forgiving.

    - **Negative:** There is an explicit or implicit clue in the text suggesting that the speaker is in a negative state, i.e., sad, angry, anxious, and violent.

    - **Neutral:** There is no explicit or implicit indicator of the speaker's emotional state: Examples are asking for like or subscription or questions about the release date or movie dialogue. This state can be considered as a neutral state.

    - **Mixed feelings:** There is an explicit or implicit clue in the text suggesting that the speaker is experiencing both positive and negative feeling: Comparing two movies

    - **Non-Malayalam:** For Malayalam if the sentence does not contain Malayalam then it is not Malayalam.

There were a total of 14 batches, where the first 13 batches had 100 items respectively and the last batch had 48 items to make a combined total of 1348 items.

116

| Metric | Original | Updated |
|---|---|---|
| Language Pair | Malayalam-English | Malayalam-English |
| Number of Tokens | 70,075 | 61,022 |
| Vocabulary Size | 19,992 | 19,389 |
| Number of Samples | 6,739 | 6,739 |
| Number of Sentences | 7,743 | 7,787 |
| Average Sentence Length | 10 | 8.26 |
| Average Sentences Per Sample | 1 | 1.15 |

Table 2: Comparison of corpus statistics reported by Chakravarthi et al. (2020) and our analysis.

| Corpus | Before | After |
|---|---|---|
| Test Size | 1,348 | 1,181 |
| Train Size | 4,851 | 4,283 |
| Validation Size | 540 | 463 |
| Total Size | 6,739 | 5,927 |

Table 3: Comparison of corpus partition sizes before and after preprocessing.

## 4.5. Preprocessing

While the preprocessing steps were outlined in the original study, it doesn't specify the packages used. Preprocessing is typically done before classifier training to prepare the data. However, according to the original study, the preprocessing phase was conducted before the data was made available to the annotators, as they illustrate in Figure 1, this was done to make annotation easier for the annotators. This motivation is unclear, as intuitively including emojis may provide more context, particularly for those examples where the sentiment is ambiguous. Data statistics are reported based on this preprocessed corpus. To confirm whether the provided data had already undergone the steps mentioned as per the original study we conducted the following preprocessing steps:

- **Removing emojis:** We removed emojis using the `emoji` package.

- **Sentence length adjustment:** We removed items with more than 15 words or less than 5 words with the `NLTK` tokeniser.

We maintain the test, train, and validation splits from the original study online.[4] However, after performing the preprocessing detailed above, the total number of samples have been reduced from 6,738 to 5,927. There were 309, 510, and 301 sentences that contained emojis, sentences exceeding 15 words and sentences fewer than 5 words, respectively. The data statistics are detailed in Table 3.

Figure 1 is taken from the original study, and shows that the preprocessing steps were conducted prior to the start of the annotation process.

This means that the provided labelled data was expected to have undergone the process of removing emojis and sentences exceeding the sentence length criteria. The descrepancies observed in this post-preprocessed data indicate that there are deviations between the actual preprocessing and preprocessing steps reported in the original study, and this in turn raises consistency issues for the data we use in the reproduction study. Hence, we decided not to perform any preprocessing steps to preserve the same corpus size before commencing the comparative corpus analysis and feeding the data to the classifiers.

## 4.6. Comparative Corpus Analysis

Comparison of the original and updated corpus statistics are detailed in Table 2.

Analysis led to some observations that are slightly different from the original findings, possibly due to variations in tools used for preprocessing and analysis. These are outlined as follows:

- **Preprocessing:** Revisiting the earlier observations, the presence of emojis and sentences exceeding specified length criteria (before making the data available) highlights the preprocessing discrepancies that we found in the original study.

- **Corpus splits:** According to the original study, the corpus includes 6,739 comments or posts. This corpus was further divided into 20% for testing (i.e., 1,348), 10% for validation (i.e., 674) and remaining 70% for training. However, upon reviewing the data provided by the original study,[5][6] we did not find this reported distribution. The data provided online has the following characteristics: while the test set contained the expected 20% of data (i.e., 1,348 items), the validation set had only 8.01% of data (i.e., 540 items), and the training set comprised 71.98% of data (i.e., 4851 items).

---

[4] https://dravidian-codemix.github.io/2020/datasets.html

[5] https://dravidian-codemix.github.io/2020/datasets.html

[6] https://github.com/bharathichezhiyan/MalayalamMixSentiment

Figure 1: Data Collection Process Of Original Study (Chakravarthi et al., 2020)

- **Label imbalance:** A notable observation regarding this corpus is its imbalanced nature, with the distribution of labels heavily skewed. Specifically, the *Positive* and *Neutral* labels are significantly over represented with 41.71% of data (i.e., 2,811 items), and 28.24% of data (i.e., 1,903 items), respectively. The *Non-Malayalam*, *Negative* and *Mixed feelings* labels have only 13.12% of data (i.e., 884 items), 10.95% of data (i.e., 738 items), and 5.98% of data (i.e., 403 items), respectively. This imbalance could have implications for the performance of the sentiment analysis classifiers that are trained on this data, as they may be biased towards the more heavily represented labels.

- **Tokenisation:** We found 61,002 tokens in contrast to the 70,075 tokens reported in the original data statistics. This variation may be due to the differences in tokenisation process followed or due to the inclusion/exclusion of specific characters as tokens. We used the word and sentence tokenisers from NLTK.[7]

  We found 7,787 sentences, while the original data statistics reported 7,743. We found a vocabulary size of 19,389 compared to 19,992 reported in the original data statistics. The average sentence length observed was 8.26. These variations may be due to the differences in tokenisation processes or to the inclusion/exclusion of specific characters as tokens.

These observations do not diminish the value of the original corpus. Rather, they highlight the complexities and challenges of working with natural language, especially in a code-mixed environment.

## 5. Classification Models

To compare the effect of re-annotation on the downstream task, we reimplement two of the original supervised classification experiments.

The classifiers we apply are:

- **Logistic Regression (LR):** The choice to utilise LR was because of its simplicity and interpretability, and because it achieved some of the best results reported by Chakravarthi et al. (2020).

- **Bidirectional Encoder Representations from Transformers (BERT)** (Devlin et al., 2019): In addtition to its good performance in the original study, we used a multilingual BERT model due to its ability to handle the multilingual aspects of the corpus. The original study fails to specify the specific BERT classifier that was used. Given this lack of detail regarding the BERT classifier, we opted for `bert-base-multilingual-uncased`[8].

## 6. Results

We report results of the annotation reproduction study in subsection 6.1, and of automated sentiment classification in subsection 6.2.

### 6.1. Human Labelling

We observed notable shifts in the label counts across all labels, as shown in Table 4. There was an increase in the label counts of labels, *Positive*, *Negative* and *Mixed feelings*, and a decrease in those of the other two, *Neutral* and *Non-Malayalam*.

| Label | Original | Re-annotated |
|---|---|---|
| *Positive* | 565 | 626 |
| *Negative* | 138 | 162 |
| *Mixed feelings* | 70 | 144 |
| *Neutral* | 398 | 327 |
| *Non-Malayalam* | 177 | 89 |

Table 4: Comparison of original and re-annotated labels for each class for the test set.

The original corpus is reported to have a Krippendorff's alpha above 0.8, indicating a high level of agreement between the annotators across the whole corpus. However, our re-annoatation of the

---

[7] https://www.nltk.org

[8] https://huggingface.co/google-bert/bert-base-multilingual-uncased

118

test corpus yielded an alpha of only 0.383. This lower score signifies that there is notable annotator disagreement within the test corpus, highlighting the challenge of achieving label consistency. This disagreement can be seen in different rounds or steps of annotation as there were 592, 175, 61 and 21 label disagreements in the annotation process steps from one till four. Although these scores are not directly comparable due to the difference in size of the test corpus and the corpus as a whole, this outcome sheds light on potential inconsistencies in annotation reliability.

## 6.2. Classification Results

**Original results**   We began by examining the performance of the `LR` and `BERT` classifiers reported in the original study. The outcomes of the original research are shown in Table 7, `BERT` achieving better preformance. The labels with the highest recall score for `LR` and `BERT` classifiers are, *Positive* and *Non-Malayalam*, respectively, suggesting its effectiveness in identifying those labels.

**Reproduction results**   For both classifiers, we evaluated their performance on the original test corpus and our re-annotated test corpus, the evaluated performance of the classifiers are shown in Table 8.

When applying the re-implemented classifiers to the original test corpus, we observed similar results to that of the classifier reported by the original study, as seen in Table 5. This indicates that the re-implementation of both classifiers can be deemed successful, and the classifiers can now be utilised to conduct a comparative analysis on both corpora.

|  | **Original** | | **Re-implemented** | |
|---|---|---|---|---|
| Classifier | LR | BERT | LR | BERT |
| macro | 0.58 | 0.61 | 0.54 | 0.65 |
| weighted | 0.66 | 0.75 | 0.63 | 0.71 |

Table 5: Results obtained by re-implementing the 2 best classifiers using the original corpus, compared to the results given in Chakravarthi et al. (2020). Note, detailed results from the original work are given in Table 7.

The comparative analysis of the re-implemented classifiers on the original and re-annotated test corpora yielded the results that are detailed in Table 8. The analysis indicates that there is a decrease in the performance of both the classifiers, as seen in Table 6.

However, the `BERT` classifier suffered a greater decrease in performance. While `LR` relies on feature engineering and does not have any multilingual understanding capabilities, `BERT` is dependent on context and subtleties within the language, and

|  | **% Decrease** | |
|---|---|---|
| Average | LR | BERT |
| macro | 11.11% | 20.00% |
| weighted | 4.76% | 14.08% |

Table 6: Average F1-score decline: Classifier results on re-annotated vs. original corpus.

might be more sensitive to modifications within the data like the labels that are assigned. This justifies the performance drop and implies that the `BERT` classifier has capabilities for capturing linguistic intricacies. This variance highlights the influence of annotation guidelines and newly annotated labels on the classifier performance.

Furthermore, given the corpora's imbalance, with the *Positive* label having the highest number of instances, it can be observed this label has the highest recall rate among all labels for both classifiers and across both corpora, indicating effectiveness of the classifiers in identifying the sentences with a positive sentiment. On the other hand, the *Mixed feelings* label exhibits the lowest recall rate across both corpora, indicating that the classifiers struggle to identify sentences with mixed sentiment.

## 6.3. Quantified Reproducibility Assessment Results

We report the reproducibility results following Belz et al. (2022). We report Type I results via coefficient of variation (CV*) and Type III results via Krippendorff's alpha ($\alpha$).

### 6.3.1. Type I Results

The comparison between the original classifier and the re-implemented classifier performance (on the original test corpus) was done using CV* (Belz, 2022; Belz et al., 2022). This was calculated using the F1-Scores of both classifiers as detailed in Table 9. Overall, the low CV* values for the macro and weighted averages of the F1-scores indicate **moderate reproducibility of the classifiers**.

Moving forward, the CV* of the re-implemented classifier performance on the original versus the re-annotated corpus was calculated, detailed in Table 10. Overall, the LR model demonstrated a CV* of 11.73 and 4.86, and the BERT model showed a CV* of 22.16, and 15.11, for the macro and weighted averages, respectively. In summary, these values suggest **less reproducibility regarding the data labels**.

### 6.3.2. Type III Results

To report Inter-Study Agreement assessment, the labels of the original test corpus and re-annotated

| | LR | | | | BERT | | | |
|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | F1-Score | Support | Precision | Recall | F1-Score | Support |
| *Mixed feelings* | 0.59 | 0.23 | 0.33 | 70 | 0.00 | 0.00 | 0.00 | 70 |
| *Negative* | 0.70 | 0.45 | 0.55 | 138 | 0.57 | 0.55 | 0.56 | 138 |
| *Neutral* | 0.65 | 0.65 | 0.65 | 398 | 0.73 | 0.79 | 0.76 | 398 |
| *Non-Malayalam* | 0.69 | 0.58 | 0.63 | 177 | 0.87 | 0.93 | 0.90 | 177 |
| *Positive* | 0.68 | 0.83 | 0.75 | 565 | 0.83 | 0.87 | 0.85 | 565 |
| macro avg | 0.66 | 0.55 | 0.58 | 1348 | 0.60 | 0.63 | 0.61 | 1348 |
| weighted avg | 0.67 | 0.67 | 0.66 | 1348 | 0.73 | 0.78 | 0.75 | 1348 |

Table 7: Results of the two best performing classifiers copied from Chakravarthi et al. (2020).

| | LR (Original Corpus) | | | | LR (Re-annotated Corpus) | | | |
|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | F1-Score | Support | Precision | Recall | F1-Score | Support |
| *Mixed feelings* | 0.80 | 0.17 | 0.28 | 70 | 0.73 | 0.08 | 0.14 | 144 |
| *Negative* | 0.77 | 0.36 | 0.49 | 138 | 0.80 | 0.31 | 0.45 | 162 |
| *Neutral* | 0.66 | 0.60 | 0.63 | 398 | 0.58 | 0.64 | 0.61 | 327 |
| *Non-Malayalam* | 0.74 | 0.51 | 0.61 | 177 | 0.39 | 0.54 | 0.45 | 89 |
| *Positive* | 0.62 | 0.85 | 0.72 | 565 | 0.68 | 0.85 | 0.76 | 626 |
| macro avg | 0.72 | 0.50 | 0.54 | 1348 | 0.64 | 0.48 | 0.48 | 1348 |
| weighted avg | 0.67 | 0.65 | 0.63 | 1348 | 0.66 | 0.63 | 0.60 | 1348 |
| | BERT (Original Corpus) | | | | BERT (Re-annotated Corpus) | | | |
| *Mixed feelings* | 0.42 | 0.44 | 0.43 | 70 | 0.42 | 0.22 | 0.29 | 144 |
| *Negative* | 0.68 | 0.51 | 0.59 | 138 | 0.69 | 0.44 | 0.54 | 162 |
| *Neutral* | 0.66 | 0.71 | 0.68 | 398 | 0.50 | 0.66 | 0.57 | 327 |
| *Non-Malayalam* | 0.81 | 0.75 | 0.78 | 177 | 0.36 | 0.65 | 0.46 | 89 |
| *Positive* | 0.77 | 0.79 | 0.78 | 565 | 0.78 | 0.72 | 0.75 | 626 |
| macro avg | 0.67 | 0.64 | 0.65 | 1348 | 0.55 | 0.54 | 0.52 | 1348 |
| weighted avg | 0.72 | 0.72 | 0.71 | 1348 | 0.63 | 0.61 | 0.61 | 1348 |

Table 8: Classifier performance on the re-annotated corpus compared to the original corpus. Note, results are reported on the test set given that our reproduction study focuses on the test set labels.

test corpus are compared by calculating Krippendorff's alpha ($\alpha$). The results are $\alpha = 0.43$. This score indicates only a moderate agreement between the original and re-annotated labels, and further suggests that there is some variability in the label consistency of the data. A detailed discussion about the label consistency is given subsequently.

## 7. Discussion

**Label Differences** The comparison of the label distribution between the original and re-annotated corpora highlight the label differences, as seen in Table 4. The labels *Mixed feelings* and *Non-Malayalam* saw significant variation, with an addition of 74 items and a reduction of 88 items, respectively. The variation in the *Mixed feelings* label implies that the instructions of how to assess sentiment complexity in the guidelines is unclear. Similarly, the discrepancy in the *Non-Malayalam* label suggests that there is a possible confusion among the annotators as to what qualifies as code-mixed and purely content that is not Malayalam. For example, the following examples are instances with high disagreement among annotators:

**Example 1:**
Tamil and Telugu padam pole aayalo...
Don't kill malayalam movies reality
*Its similar to Tamil and Telugu films...*
*Don't kill malayalam movies reality*
Assigned Label: *Mixed feelings*

**Example 2:**
Numma or nummade or nammande palakkad le Katha aanu
*Our own palakkads story*
Assigned Label: *Neutral*

In **example 1**, there are three different perspectives. Firstly, the sentence could be seen as a **neutral** observation where Malayalam films are being compared to Tamil and Telugu films. Secondly, the advice '*Don't kill malayalam movies reality*' implies a **negative** sentiment towards the Tamil and Telugu industries. Thirdly, the sentence might imply a positive view towards Tamil and Telugu cinema's handling of reality and then warn against the destruction of reality in Malayalam films, suggesting a **mixed sentiment**.

In **example 2**, the underlying sentiments are

| Labels | LR | | | BERT | | |
|---|---|---|---|---|---|---|
| | Original | Re-implemented | **CV\*** | Original | Re-implemented | **CV\*** |
| *Positive* | 0.75 | 0.72 | 4.07 | 0.85 | 0.78 | 8.56 |
| *Negative* | 0.55 | 0.49 | 11.50 | 0.56 | 0.59 | **5.20** |
| *Mixed feelings* | 0.33 | 0.28 | **16.34** | 0.00 | 0.43 | **199.40** |
| *Neutral* | 0.65 | 0.63 | **3.12** | 0.76 | 0.68 | 11.08 |
| *Non-Malayalam* | 0.63 | 0.61 | 3.22 | 0.90 | 0.78 | 14.25 |
| macro avg | 0.58 | 0.54 | 7.12 | 0.61 | 0.65 | 6.33 |
| weighted avg | 0.66 | 0.63 | 4.64 | 0.75 | 0.71 | 5.46 |

Table 9: Quantitative Reproducibility Analysis utilising CV\* between the results reported in the original paper and our re-implemented classifiers (using the original corpus). CV\* is calculated based on the F1-scores.

| Labels | LR | | | BERT | | |
|---|---|---|---|---|---|---|
| | Original | Re-annotated | **CV\*** | Original | Re-annotated | **CV\*** |
| *Positive* | 0.72 | 0.76 | 5.39 | 0.78 | 0.75 | **3.90** |
| *Negative* | 0.49 | 0.45 | 8.49 | 0.59 | 0.54 | 8.82 |
| *Mixed feelings* | 0.28 | 0.14 | **66.47** | 0.43 | 0.29 | 38.78 |
| *Neutral* | 0.63 | 0.61 | **3.22** | 0.68 | 0.57 | 17.55 |
| *Non-Malayalam* | 0.61 | 0.45 | 30.10 | 0.78 | 0.46 | **51.46** |
| macro avg | 0.54 | 0.48 | 11.73 | 0.65 | 0.52 | 22.16 |
| weighted avg | 0.63 | 0.60 | 4.86 | 0.71 | 0.61 | 15.11 |

Table 10: Quantitative Reproducibility Analysis utilising CV\* of the re-implemented classifiers on the original versus the re-annotated corpus (i.e. detailed in Table 8). CV\* is calculated based on the F1-scores.

| Labels | Example 1 | Example 2 |
|---|---|---|
| *Positive* | 0 | 2 |
| *Negative* | 2 | 0 |
| *Neutral* | 2 | 3 |
| *Mixed feelings* | 3 | 0 |
| *Non-Malayalam* | 0 | 2 |

Table 11: Comparison of labels assigned to example items.

*positive*, *neutral* and *non-Malayalam*. **Positive** because the phrase suggests pride to be part of the Palakkad district. Without context, the sentence could be seen as simply stating a fact, thus implying the **neutral** sentiment. Lastly, the code-mixed text can be interpreted as both Malayalam or Kannada, as '*Numma*' or '*Nummade*' are both words that are present in both languagues, this confusion can lead annotators to opt for the **non-malayalam** sentiment.

Moreover, the removal of emojis before annotation could have a significant effect on the underlying sentiment. Additionally, the challenges in code-mixed data such as the ambiguity outlined in the examples earlier could have been lessened with the help of more clear and detailed annotation guidelines.

**Issues Affecting Reproducbility**   In the process of attempting to reproduce the results of another study, we faced several significant challenges that underscore the complexities of research reproduciblity. The following list outlines the reproduction challenges that were encountered:

- **Data Preparation Issues:** Chakravarthi et al. (2020) explain that *preprocessing* efforts were conducted to alleviate potential challenges for the annotators. However, the labelled data had numerous instances that appear not to have undergone preprocessing. The discrepancy between the documentation and the provided data poses a significant challenge to the reproducibility and hinders the integrity of the preprocessed data. Moreover, the study followed a structured approach to the *annotation process* which involved a three-step process. However, in this methodology, there is a critical ambiguity in addressing scenarios where the annotators continued in disagreement beyond the third step. Unlike the study's decision to discard such data, this reproducibility challenge was addressed by taking the decision to involve a sixth annotator to resolve those disagreements, and any other pending disagreements afterwards were resolved by me. Lastly, the absence of the actual *annotation guidelines*, apart from the basic schema, presented a significant challenge. Without these guidelines annotators faced ambiguity and had varied interpretations for the same sentences.

- **Classification Issues:** Although the total size

of the provided corpus is accurate, the specified partition counts mentioned in the original study for training and validation is incorrect. This creates confusion and inconsistency in understanding the *corpus partitions*, which affects the reliability and reproducibility of the study and its corpus. Additionally, the original study asserts the free *availability of the code* and corpus for research purposes. However, this assertion is not met as the GitHub repository only contains a readme file with the corpus links but lacks the actual code. This situation complicates the replication process, stressing the importance of resource sharing in the NLP community. Moreover, acheiving comparable classifier performance given the lack of access to the original code, also posed a significant reproduciblity challenge. Furthermore, there is an uncertainity in the *classifier variant selection* for BERT in the original study. This oversight in not specifying the version was resolved by opting for the BERT-uncased-multilingual version. However, the differences in classifier version can hinder the results, thereby, affecting the reproducibility of the original study.

## 8.  Conclusion

Our findings contribute to the ongoing discussion on the reproducibility and authenticity of research conducted in the field of NLP. The reproduction study yielded results that demonstrate a decrease in the performance accuracy of the re-implemented classifiers when compared to the results of the original study. Subsequently, **we were not able to reproduce the original study's results**. The obstacles we faced were **preprocessing inconsistencies**, **lack of guidelines and code**, **unclear annotation processes**, and **missing information** regarding packages and classifier variants used in the original study.

To advance the field and mitigate these reproducbility challenges, future work should focus on the development and adoption of reporting frameworks that are standardised. Additionally, the sharing of code, corpora, and detailed methodologies should be encouraged in the NLP community and studies assessing reproduciblity should be conducted systematically to pave the way for reliable and authentic researches.

## Limitations and Ethical Considerations

**Limitations**  Although this research provides insights into the reproducibility of NLP label collection, it has several limitations. The scope of this study is limited to the test corpus and set of pre-selected ML classifiers of the original research,

which may not fully capture the underlying reproducbility challenges. Moreover, although the `bert-base-multilingual-uncased` classifier is designed to handle multiple languages, the study's approach, including the use of `LR`, neither the original nor this study explicitly addresses the intricacies of code-mixing. Furthermore, even though we were successful in re-implementing the classifiers in this study, it still might not mirror exactly those used in the original study, influencing the performance comparison and assessment of reproducbility.

**Ethical Considerations**  This study was conducted with the approval of the institutional review board of Heriot-Watt University.  Data was collected and stored on the Heriot-Watt-approved MS OneDrive system and complies with the General Data Protection Regulation (GDPR). Participant consent was obtained through an online information sheet and consent form prior to any data collection.

## Acknowledgements

## Bibliographical References

Gavin Abercrombie, Dirk Hovy, and Vinodkumar Prabhakaran. 2023. Temporal and second language influence on intra-annotator agreement and stability in hate speech labelling. In *Proceedings of the 17th Linguistic Annotation Workshop (LAW-XVII)*, Proceedings of the Annual Meeting of the Association for Computational Linguistics, pages 96–103. Association for Computational Linguistics.

Aryaman Arora, Adam Farris, Samopriya Basu, and Suresh Kolichala. 2022. Computational historical linguistics and language diversity in South Asia. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1396–1409, Dublin, Ireland. Association for Computational Linguistics.

Monya Baker. 2016. 1,500 scientists lift the lid on reproducibility. *Nature*, 533(7604):452–454.

Anya Belz. 2022. A metrological perspective on reproducibility in NLP*. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Shubham Agarwal, Anastasia Shimorina, and Ehud Reiter. 2021. A systematic review of reproducibility research in natural language processing. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 381–393, Online. Association for Computational Linguistics.

Anya Belz, Maja Popovic, and Simon Mille. 2022. Quantified reproducibility assessment of NLP results. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16–28, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *Proceedings of the Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Bharathi Raja Chakravarthi, Navya Jose, Shardul Suryawanshi, Elizabeth Sherly, and John Philip McCrae. 2020. A sentiment analysis dataset for code-mixed Malayalam-English. In *Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL)*, pages 177–184, Marseille, France. European Language Resources association.

Jacob Devlin, 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. Association for Computational Linguistics.

Penelope Gardner-Chloros. 2009. *Code-switching*. Cambridge university press.

Odd Erik Gundersen and Sigbjørn Kjensmo. 2018. State of the art: Reproducibility in artificial intelligence. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1).

H.H. Hock and E. Bashir. 2016. *The Languages and Linguistics of South Asia: A Comprehensive Guide*. The World of Linguistics [WOL]. De Gruyter.

Douwe Kiela, Changhan Wang, and Kyunghyun Cho. 2018. Dynamic meta-embeddings for improved sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1466–1477, Brussels, Belgium. Association for Computational Linguistics.

Gina-Anne Levow, Valerie Freeman, Alena Hrynkevich, Mari Ostendorf, Richard Wright, Julian Chan, Yi Luan, and Trang Tran. 2015. Recognition of stance strength and polarity in spontaneous speech. pages 236–241.

Bing Liu. 2012. *Sentiment analysis and opinion mining Bing Liu.* Synthesis digital library of engineering and computer science. Morgan & Claypool, San Rafael, Calif. (1537 Fourth Street, San Rafael, CA 94901 USA).

Saif Mohammad. 2016. A practical guide to sentiment annotation: Challenges and solutions. In *Proceedings of the 7th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 174–179, San Diego, California. Association for Computational Linguistics.

Andrew Moore and Paul Rayson. 2018. Bringing replication and reproduction together with generalisability in NLP: Three reproduction studies for target dependent sentiment analysis. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1132–1144, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

123

Mika V. Mäntylä, Daniel Graziotin, and Miikka Kuutila. 2018. The evolution of sentiment analysis—a review of research topics, venues, and top cited papers. *Computer Science Review*, 27:16–32.

Deepu S. Nair, Jisha P. Jayan, Rajeev R R, and Elizabeth Sherly. 2014. Sentima - sentiment extraction for malayalam. *2014 International Conference on Advances in Computing, Communications and Informatics (ICACCI)*, pages 1719–1723.

Nicolas P Rougier, Konrad Hinsen, Frédéric Alexandre, Thomas Arildsen, Lorena A Barba, Fabien CY Benureau, C Titus Brown, Pierre De Buyl, Ozan Caglayan, Andrew P Davison, et al. 2017. Sustainable computational science: the rescience initiative. *PeerJ Computer Science*, 3:e142.

S. Soumya and K.V. Pramod. 2020. Sentiment analysis of Malayalam tweets using machine learning techniques. *ICT Express*, 6(4):300–305.

Craig Thomson, Ehud Reiter, and Anya Belz. 2024. Common Flaws in Running Human Evaluation Experiments in NLP. *Computational Linguistics*, pages 1–11.

Janyce M. Wiebe. 1990. Identifying subjective characters in narrative. In *COLING 1990 Volume 2: Papers presented to the 13th International Conference on Computational Linguistics*.

Martijn Wieling, Josine Rawee, and Gertjan van Noord. 2018. Squib: Reproducibility in computational linguistics: Are we willing to share? *Computational Linguistics*, 44(4):641–649.

Fanghua Ye, Jarana Manotumruksa, and Emine Yilmaz. 2022. MultiWOZ 2.4: A multi-domain task-oriented dialogue dataset with essential annotation corrections to improve state tracking evaluation. In *Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 351–360, Edinburgh, UK. Association for Computational Linguistics.

Peng Zhou, Zhenyu Qi, Suncong Zheng, Jiaming Xu, Hongyun Bao, and Bo Xu. 2016. Text classification improved by integrating bidirectional LSTM with two-dimensional max pooling. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 3485–3495, Osaka, Japan. The COLING 2016 Organizing Committee.

# Reproducing the Metric-Based Evaluation of a Set of Controllable Text Generation Techniques

**Michela Lorandi, Anya Belz**

ADAPT Centre, Dublin City University, Ireland

{michela.lorandi, anya.belz}@adaptcentre.ie

## Abstract

Rerunning a metric-based evaluation should be more straightforward, and results should be closer, than in a human-based evaluation, especially where code and model checkpoints are made available by the original authors. As this report of our efforts to rerun a metric-based evaluation of a set of single-attribute and multiple-attribute controllable text generation (CTG) techniques shows however, such reruns of evaluations do not always produce results that are the same as the original results, and can reveal errors in the reporting of the original work.

**Keywords:** Reproduction, metric evaluation, controllable text generation.

## 1. Introduction

Over the past few years, the fields of natural language processing (NLP) and machine learning (ML) have seen an increase in interest in reproducibility (Sinha et al., 2020; Branco et al., 2020; Belz et al., 2021; Belz and Thomson, 2023). Initially, efforts focussed on promoting and encouraging sharing of all resources needed to rerun experiments, but increasingly it became clear that exact reproduction of results is rarely the outcome even where metric evaluation is concerned. The question is what can be concluded in such situations beyond binary reproduced vs. not reproduced findings.

Belz et al. (2022; 2023) proposed QRA++, an approach to measuring how close results from two evaluations are, and how reproducible evaluation measures are, in order to facilitate comparison in terms of degree of reproducibility between different methods of evaluation. This approach enables comparable, quantified reproducibility results to be produced.

In this short report, we present our work rerunning the metric-based evaluation of a set of single and multiple-attribute controllable text generation techniques (Gu et al., 2022, 2023). In the case of all except one pair of scores from the original and reproduction evaluations, the two scores are not the same, and we apply QRA++ to quantify the differences.

We start with a summary of the QRA++ measures we use (Section 2), followed by a description of the specific original experiments we repeated in this reproduction study (Section 3). We then describe how we went about repeating the work (Section 4), before presenting the side-by-side results from the original work and our reproduction along with the QRA++ measures of their similarity (Section 5). We finish with some discussion and conclusions (Section 6).

## 2. QRA++ Measures

QRA++ distinguishes four types of results commonly reported in NLP and ML papers:

1. Type I results: single numerical scores, e.g. mean quality rating, error count, etc.
2. Type II results: sets of related numerical scores, e.g. set of Type I results .
3. Type III results: categorical labels attached to text spans of any length.
4. Type IV results: Qualitative findings stated explicitly or implied by quantitative results in the original paper.

The above are quantitatively assessed as follows:

1. Type I results: Small-sample coefficient of variation CV* (Belz, 2022).
2. Type II results: Pearson's r, Spearman's $\rho$.
3. Type III results: Multi-rater: Fleiss's $\kappa$; Multi-rater, multi-label: Krippendorff's $\alpha$.
4. Type IV results: Proportion of findings that are / are not confirmed by the repeat experiment. To obtain comparable results we restrict ourselves to pairwise system ranks as findings.

In the work reported in this paper we have Type I, II and IV results, and therefore apply the corresponding quantitative measures above.

## 3. Original Work Being Repeated

In the present reproduction study, we carried out repeat evaluations of the main new systems presented by Gu et al. (2022) and Gu et al. (2023). The authors provide the code on GitHub[1] and the model checkpoints on Google Drive.[2]

---

[1] https://github.com/HappyGu0524/Multi Control

[2] https://drive.google.com/drive/folde rs/14XHSG4IAGlAL9t-SYoTUKnAs5ARqHd5f

More precisely, the experimental grid we reproduced looks as follows: {PriorCTG} x {Topic (World, Sports, Business, and Technology), Sentiment (Positive and Negative), Toxicity (Toxic and Non-Toxic)} x {no extension, extension} + {PriorCTG} x {Multi attribute (Topic, Sentiment and Non-Toxic)} x {no optim, optim} + {MultiCTG} x {Multi attribute (Topic, Sentiment and Non-Toxic)}. The individual systems (MultiCTG, PriorCTG +/- extend/optim) are described in the next section.

## 3.1. Systems included in reproduction

We included the results for the four main new systems from the original work (Gu et al., 2022, 2023) in our reproduction study; we abbreviate system names as follows: MultiCTG, PriorCTG, PriorCTG+extend, and PriorCTG+optim.

*MultiCTG*: This is the core new CTG approach proposed by Gu et al. (2022) which directly searches for the intersection areas of multiple attribute distributions to achieve control over multiple control attributes. The attribute space is first estimated with an autoencoder structure, then the intersections are iteratively approached via joint minimisation of distances to points representing the controlled attributes.

*PriorCTG*: This is the core new CTG approach proposed by Gu et al. (2023), which utilises a form of latent-space control, more specifically an invertible transformation function, the Normalizing Flow, that maps the complex distributions in latent space to simple Gaussian distributions in **prior** space.

*PriorCTG+extend*: The **extend** control strategy additionally achieves opposite control, as in contrastive learning, by using negative weights when interpolating.

*PriorCTG+optim*: The **optim** control strategy additionally optimises the intersection of the single attribute representations in prior space to achieve multiple-attribute control.

All systems are trained on the IMDb movie reviews dataset (Maas et al., 2011), the AGNews dataset (Zhang et al., 2015), and the Jigsaw Toxic Comment Classification Challenge Dataset (cjadams, 2017), respectively, for control of sentiment, topic and detoxification attributes. Note that we did not include any of the baseline systems in the reproduction.

## 3.2. Evaluation metrics

The metrics in this section are all described in detail in Gu et al. (2022). The main set of metrics assesses single-attribute control performance (called 'attribute relevance' in the original papers), computed as the percentage of outputs that are classified as having the given intended control attribute value by a specific classifier.

For *Sentiment* control performance, the classifier is DeBERTa (He et al., 2020) finetuned on the Yelp dataset (Zhang et al., 2015).

For *Topic* control performance, the classifier is DeBERTa finetuned on the AGNews dataset (Zhang et al., 2015) utilizing the portion of dataset not used during the model's training.

For *Toxicity* control performance, there is a discrepancy between what the paper says and what is in the evaluation script shared on GitHub. According to the former, toxicity is measured with the Google Perspective API.[3] However, the script uses a toxicity classifier obtained by finetuning DeBERTa on the Jigsaw Toxic Comment Classification Challenge Dataset,[4] analogous to control performance assessment for the other control attributes. We ran the evaluation both with Perspective and with the DeBERTa classifier, and found that scores obtained with the latter were closer to the original scores, so those are what we used.

*Multiple-attribute control performance* is computed as the average of the single-attribute control performance scores for the three attributes being controlled.

*Perplexity* is calculated by GPT2-large following the Contrastive Prefix method (Qian et al., 2022). Note that we used our own implementation as no code was shared for this.

*Distinctness* (Li et al., 2016) is computed as the percentage of distinct n-grams in the continuations generated from a given set of prefixes. System-level 1-gram, 2-gram, and 3-gram distinctness scores are obtained by averaging over prefix-level distinctness scores. In multi-control setting, the average of system-level Distinct-1, 2 and 3 is computed. Here too we used our own implementation based on Yu et al. (2021) implementation, because the code was not shared either by Li et al. or by Gu et al.

This gives us six main types of metrics (the three classifier-based metrics, their average (for multiple-attribute control), perplexity, and distinctness). In Table 1 we additionally give the average over the individual control performance scores (**Avg.** columns) for sentiment, topic and toxicity.

## 4. Reproduction Work

Our first step was to download the code and model checkpoints from the authors' Github and Drive repositories, and recreate the environments on our machine with a GPU RTXA6000 with 48GB RAM.

We then re-executed the inference phase of the experiments involving PriorCTG from Gu et al. (2023), first those with single-attribute control, i.e.

---

[3] https://www.perspectiveapi.com/
[4] https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/

| Methods | Sentiment↑ (%) | | | Topic↑ (%) | | | | | Detox.↑ (%) | PPL.↓ | Dist.-1/2/3↑ (%) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Avg. | Pos. | Neg. | Avg. | W. | S. | B. | T. | | | |
| PriorCTG | 97.1 | 99.9 | 94.3 | 95.9 | 95.5 | 99.3 | 90.2 | 98.7 | 90.7 | 61 | 42.0 / 79.7 / 88.4 |
| PriorCTG Repro | 98.2 | 99.9 | 96.6 | 94.8 | 93.4 | 97.8 | 88.5 | 99.5 | 96.9 | 59.7 | 41.9 / 79.5 / 88.4 |
| PriorCTG+extend | 99.7 | 99.9 | 99.5 | 97.8 | 97.9 | 99.4 | 94.0 | 99.8 | 95.7 | 61.6 | 42.4 / 79.4 / 88.1 |
| PriorCTG+extend Repro | 99.3 | 99.9 | 98.7 | 98.2 | 98.2 | 99.5 | 95.5 | 99.8 | 99.9 | 60.8 | 42.3 / 79.2 / 88.1 |

Table 1: Side-by-side metric results from original work (Gu et al., 2023) and reproduction study for **single-attribute control** (last two rows in Table 1 in the original paper). The results of the last two columns are obtained using our own implementation. For PriorCTG and PriorCTG+extend systems (see Section 3). Repro=Reproduction results.

| Methods | Average↑ (%) | Sentiment↑ (%) | Topic↑ (%) | Detoxification↑ (%) | PPL.↓ | Dist.↑ (%) |
|---|---|---|---|---|---|---|
| MultiCTG | 87.4 ± 10.9 | 86.7 ± 10.5 | 84.8 ± 14.2 | 90.7 ± 7.4 | 31.3 | 59.0 |
| MultiCTG Repro | 88.4 ± 8.3 | 84.9 ± 11.5 | 84.5 ± 14.4 | 95.9 ± 5.5 | 31.5 | 59.2 |
| PriorCTG | 89.9 ± 8.7 | 88.0 ± 10.6 | 87.4 ± 8.5 | 94.3 ± 3.2 | 38.9 | 65.3 |
| PriorCTG Repro | 91.1 ± 6.7 | 88.0 ± 10.2 | 87.1 ± 11.2 | 98.3 ± 1.6 | 38.3 | 65.2 |
| PriorCTG+optim | 92.2 ± 8.6 | 92.5 ± 8.5 | 89.3 ± 11.0 | 94.9 ± 3.4 | 33.0 | 61.7 |
| PriorCTG+optim Repro | 93.2 ± 7.2 | 91.8 ± 9.7 | 89.3 ± 12.4 | 98.6 ± 1.1 | 32.5 | 62 |

Table 2: Side-by-side metric results from original work (Gu et al., 2022, 2023) and reproduction study for **multiple-attribute control**. Results for MultiCTG are from the third to last row in Gu et al. (2022). Original results for the other two systems are from the last two rows in Table 3 in Gu et al. (2023). The results of the last two columns are obtained using our own implementation. For system and metrics descriptions see Section 3). Repro=Reproduction results.

where Topic, Sentiment or Toxicity are being controlled individually, and then those with multiple-attribute control, where Topic, Sentiment and Toxicity are being controlled at the same time.

For multiple-attribute control we also re-executed the inference phase of the experiments involving MultiCTG from Gu et al. (2022). This gave us sets of $35 \times 5 = 175$ outputs (35 inputs from the PPLM Prompts test set $\times$ 5 repetitions of prompting and collecting the outputs) for each system/attribute combination.

Note that as in the original work, outputs are generated for all values of all controlled attributes (single-attribute case) or for all combinations of controlled attribute values (multiple-attribute case), results for all of which except Toxicity=toxic ('Detox(ification)' in the tables) are reported in the results tables. In the multiple-attribute case, the average over different attribute value combinations, along with the corresponding standard deviation, is reported.

For the evaluation, we computed the metrics listed in Section 3. Recall from Section 3.2 that we used the script provided by the authors for Sentiment, Topic and Toxicity control performance assessment. However, we coded our own scripts to compute Perplexity and Distinct-n, as scripts are not provided for these. We also use our own code for the standard deviations in the multiple-attribute table. For all scripts we use parameters as provided by the authors.

Note that as a result of some of the evaluation scripts not being shared, we have two distint reproduction situations (which in QRA++ is reflected in the measurement conditions): (a) for the classifier-based control-performance measures, we use our outputs (regenerated by us using the original authors' code) and evaluate them with the original authors' scripts; and (b) for perplexity and distinctness, we use our outputs *and* our evaluation scripts. In the former case differences in scores can only be due to differences in *executing* the original authors' code, whereas in the latter case, differences can be due to both execution and differences in the evaluation code.

In order to avoid this dual possible source of difference for perplexity and distinctness scores, we decided to re-evaluate the original authors' outputs with our own script. This means that the scores in our tables are not the same as in the two original papers for these two metrics. But it means CV* scores and other reproducibility measures are comparable across all metrics.

## 5. Side-by-Side Results and QRA++ Assessment

Tables 1 and 2 present side-by-side evaluation results for the original and reproduction work, for each of the six metrics from Section 3, plus, in Table 1 only, averages over individual control performance scores (**Avg.** columns). Reall that we reevaluated the original authors' outputs in terms of Perplexity and Distinctness (see preceding section).

| System | CV* between original and reproduction scores for each evaluation measure | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Sent avg | Sent pos | Sent neg | Topic avg | Topic W | Topic S | Topic B | Topic T | Detox | PPL | Dist-1 | Dist-2 | Dist-3 |
| Prior-CTG | 1.12 | 0 | 2.4 | 1.15 | 2.22 | 1.52 | 1.9 | 0.8 | 6.59 | 2.15 | 0.24 | 0.25 | 0 |
| Prior-CTG+ext | 0.4 | 0 | 0.8 | 0.41 | 0.31 | 0.1 | 1.58 | 0 | 4.28 | 1.3 | 0.24 | 0.25 | 0 |
| Average | 0.76 | 0 | 1.6 | 0.78 | 1.27 | 0.81 | 1.74 | 0.4 | 5.44 | 1.725 | 0.24 | 0.25 | 0 |

Table 3: CV* for each pair of original and reproduction metric scores, for the Prior-CTG and Prior-CTG+extend systems, and the average over both systems.

| System | CV* between original and reproduction scores for each evaluation measure | | | | | |
|---|---|---|---|---|---|---|
| | Avg | Sentiment | Topic | Detox | PPL | Distinct-n |
| Multi-CTG | 1.13 | 2.09 | 0.35 | 5.56 | 0.64 | 0.34 |
| Prior-CTG | 1.32 | 0.0 | 0.34 | 4.14 | 1.52 | 0.15 |
| Prior-CTG+optim | 1.08 | 0.76 | 0.0 | 3.81 | 1.52 | 0.48 |
| Average | 1.18 | 0.95 | 0.23 | 4.5 | 1.23 | 0.32 |

Table 4: CV* for each pair of original and reproduction metric scores, for the Multi-CTG, Prior-CTG and Prior-CTG+optim systems, and the average over all three.

### 5.1. Type IV results

Regarding Type IV results (findings), here we are assessing relative performance between systems, such that each pairwise ranking counts as one finding. Note that statistical significance was not computed in the original work.

For single-attribute control (Table 1), in the original work, Prior CTG+extend has higher scores than PriorCTG according to all metrics except for Perplexity and 2-gram and 3-gram Distinctness where PriorCTG scores are very slightly higher. For Sentiment/Pos, scores are identical. In our reproduction evaluations, these two systems are ranked the same way in all cases, giving us a perfect proportion of 13/13 findings upheld for this table.

For multiple-attribute control (scores in Table 2), the same type of analysis gives us a proportion of 18/18 findings upheld (pairwise ranks confirmed).

### 5.2. Type I results

For Type I results, we computed CV* values for all individual system/metric level original and reproduction scores. We report the individual scores, as well as the mean per metric.

For single-attribute control (scores in Table 1), Table 3 shows CV* scores for each pair of original and reproduction metric scores, for the Prior-CTG and Prior-CTG+extend systems, and the average over both systems (last row).

One clear tendency is that the Prior-CTG system has better reproducibility scores across the board than Prior-CTG+extend (except for distinctness metrics where the two systems are tied).

Looking at metric-level differences ('Average' row), we can see that Perplexity and (by a smaller margin) Detoxification Control have lower reproducibility than the other metrics.

For multiple-attribute control (scores in Table 2), Table 4 shows CV* scores for each pair of original and reproduction metric scores, for the Multi-CTG, Prior-CTG and Prior-CTG+optim systems, and the average over all three (last row). We can see that here too, the Perplexity and Detoxification Control metrics have the poorest reproducibility.

We can also see a slight tendency for the classifier scores for the Prior-CTG+optim system to have better reproducibility than the other two systems (but not for PPL and Distinct-n), but the picture is more mixed than for the single-attribute control systems.

### 5.3. Type II results

For Type II results we compute Pearson's correlation coefficients between sets of metric scores in two ways, (i) for each metric (i.e. how do all the scores for each metric correlate between original and reproduction), and (ii) for each system (i.e. how do all the scores for each system correlate).

For single-attribute control (scores in Table 1), system-level Pearson's between all metric results in the original and reproduction runs is above 0.99 for both Prior-CTG and Prior-CTG+extend. Mean metric-level Pearson's is perfect (but note that we have only two score pairs all of which are ranked identically).

For multiple-attribute control (scores in Table 2), system-level Pearson's between all metric results in the original and reproduction runs is above 0.99 for all three systems. Metric-level Pearson's is above 0.99 for all metrics except the sentiment-classifier metric which at $r = 0.969$ is slightly lower than the

other metrics. Mean metric-level $r$ is 0.994.

## 6. Discussion and Conclusion

The main challenges in carrying out our reproduction study were (i) lack of clarity in the paper with respect to what the averages and standard deviations in results tables were computed over, and (ii) discrepancies between the shared code and what the paper said, e.g. the paper says toxicity was assessed with Perspective, whereas the shared evaluation script has a toxicity classifier.

Our quantified reproducibility assessments revealed a high degree of reproducibility at the study level for Type II and Type IV results. For Type I results, study-level CV* (computed as the mean of metric-level means) was 1.154 for single-attribute control, and 1.402 for multiple-attribute control. While this compares well to reproducibility results in human evaluations which very rarely achieve study-level CV* below 5 in pairwise comparisons of original study and one reproduction, it does confirm once again that even with identical code, we cannot necessarily expect to get the same results.

In terms of metric-level CV*, the Detoxification control metric had notably worse reproducibility than the others which may be partly but not entirely explainable by the fact that only Toxicity=nontoxic was taken into account here.

In terms of the results that tend to be considered as most important, Type IV results or findings upheld, reproducibility was perfect with all pairwise rankings being identical in the original and reproduction experiments.

## Bibliographical References

Anya Belz. 2022. A metrological perspective on reproducibility in nlp. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Shubham Agarwal, Anastasia Shimorina, and Ehud Reiter. 2021. A systematic review of reproducibility research in natural language processing. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 381–393.

Anya Belz and Craig Thomson. 2023. The 2023 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 35–48.

António Branco, Nicoletta Calzolari, Piek Vossen, Gertjan Van Noord, Dieter Van Uytvanck, Joao Silva, Luís Gomes, André Moreira, and Willem Elbers. 2020. A shared task of a new, collaborative type to foster reproducibility: A first exercise in the area of language science and technology with reprolang2020. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 5539–5545. European Language Resources Association (ELRA).

Julia Elliott Lucas Dixon Mark McDonald nithum Will Cukierski cjadams, Jeffrey Sorensen. 2017. Toxic comment classification challenge.

Yuxuan Gu, Xiaocheng Feng, Sicheng Ma, Lingyuan Zhang, Heng Gong, and Bing Qin. 2022. A distributional lens for multi-aspect controllable text generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1023–1043, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Yuxuan Gu, Xiaocheng Feng, Sicheng Ma, Lingyuan Zhang, Heng Gong, Weihong Zhong, and Bing Qin. 2023. Controllable text generation via probability density estimation in the latent space. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12590–12616, Toronto, Canada. Association for Computational Linguistics.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*.

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San Diego, California. Association for Computational Linguistics.

Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.

Jing Qian, Li Dong, Yelong Shen, Furu Wei, and Weizhu Chen. 2022. Controllable natural language generation with contrastive prefixes. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2912–2924,

Dublin, Ireland. Association for Computational Linguistics.

Koustuv Sinha, Joelle Pineau, Jessica Forde, Rosemary Nan Ke, and Hugo Larochelle. 2020. Neurips 2019 reproducibility challenge. 6.

Kirstie Whitaker. 2017. The MT Reproducibility Checklist. https://www.cs.mcgill.ca/~j pineau/ReproducibilityChecklist.pd f.

Dian Yu, Zhou Yu, and Kenji Sagae. 2021. Attribute alignment: Controlling text generation from pretrained language models. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2251–2268, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.

## A. Perplexity and Distinct-n implementation

No code was shared to compute perplexity and Distinct-n, hence we used our own implementation. Perplexity is calculated using the evaluate library of HuggingFace[5] using GPT-2 Large.

System-level Distinct-n (n=1, 2, 3) is the average Distinct-n at prefix-level, which is computed as the number of unique n-grams in the set of generated outputs with the same prefix over the total amount of tokens. GPT-2 is used to tokenise the texts.

Table 5 and 6 show Perplexity and Distinct-n results reported in the original work, the results of the original study computed using our implementation and the reproduction using our implementation.

---

[5] https://huggingface.co/docs/evaluate /en/index

| Methods | PPL. (%) ↓ | Dist.-1/2/3↑ (%) |
|---|---|---|
| PriorCTG | 54.3 | 29.1 / 70.1 / 86.9 |
| PriorCTG using our implementation | 61 | 42.0 / 79.7 / 88.4 |
| PriorCTG Repro | 59.7 | 41.9 / 79.5 / 88.4 |
| PriorCTG+extend | 54.6 | 29.8 / 70.5 / 86.8 |
| PriorCTG+extend using our implementation | 61.6 | 42.4 / 79.4 / 88.1 |
| PriorCTG+extend Repro | 60.8 | 42.3 / 79.2 / 88.1 |

Table 5: Side-by-side metric results from original work (Gu et al., 2023), original work (Gu et al., 2023) computed using our own implementation and reproduction study using our own implementation for **single-attribute control** (last two rows in Table 1 in the original paper). For PriorCTG and PriorCTG+extend systems (see Section 3). Repro=Reproduction results.

| Methods | PPL.↓ | Dist.↑ (%) |
|---|---|---|
| MultiCTG | 28.4 | 49.5 |
| MultiCTG using our implementation | 31.3 | 59.0 |
| MultiCTG Repro | 31.5 | 59.2 |
| PriorCTG | 34.7 | 55.5 |
| PriorCTG using our implementation | 38.9 | 65.3 |
| PriorCTG Repro | 38.3 | 65.2 |
| PriorCTG+optim | 29.6 | 51.6 |
| PriorCTG+optim using our implementation | 33.0 | 61.7 |
| PriorCTG+optim Repro | 32.5 | 62 |

Table 6: Side-by-side metric results from original work (Gu et al., 2022, 2023), original work (Gu et al., 2022, 2023) computed using our own implementation and reproduction study using our own implementation for **multiple-attribute control**. Results for MultiCTG are from the third to last row in Gu et al. (2022). Original results for the other two systems are from the last two rows in Table 3 in Gu et al. (2023). For system and metrics descriptions see Section 3). Repro=Reproduction results.

# ReproHum: #0033-03: How Reproducible Are Fluency Ratings of Generated Text? A Reproduction of August et al. 2022

**Emiel van Miltenburg, Anouck Braggaar, Nadine Braun, Martijn Goudbeek,**
**Emiel Krahmer, Chris van der Lee, Steffen Pauws, Frédéric Tomas**

Tilburg center for Cognition and Communication, Tilburg University
Warandelaan 2, 5037 AB Tilburg, The Netherlands
✉ C.W.J.vanMiltenburg@tilburguniversity.edu

## Abstract

In earlier work, August et al. (2022) evaluated three different Natural Language Generation systems on their ability to generate *fluent*, *relevant*, and *factual* scientific definitions. As part of the ReproHum project (Belz et al., 2023), we carried out a partial reproduction study of their human evaluation procedure, focusing on human fluency ratings. Following the standardised ReproHum procedure, our reproduction study follows the original study as closely as possible, with two raters providing 300 ratings each. In addition to this, we carried out a second study where we collected ratings from eight additional raters and analysed the variability of the ratings. We successfully reproduced the inferential statistics from the original study (i.e. the same hypotheses were supported), albeit with a lower inter-annotator agreement. The remainder of our paper shows significant variation between different raters, raising questions about what it really means to reproduce human evaluation studies.

**Keywords:** Fluency ratings, Natural Language Generation, Evaluation, Reproduction

## 1. Introduction

The quality of automatically generated texts is often evaluated using human ratings because they allow us to assess a wide range of different kinds of quality dimensions, ranging from FORM (*grammaticality, fluency*) to CONTENT (*correctness, appropriateness*) and SUITABILITY FOR PARTICULAR PURPOSES (*usability, informativeness*). One major challenge for Natural Language Generation (NLG) research is to properly define and operationalise all the different constructs that one may be interested in (Belz et al., 2020). At the moment, there is a lack of standardisation in the field, leading to terminological confusion. *Fluency* is a major culprit; Howcroft et al. (2020) show how different authors use the same term to refer to *fifteen* different constructs. Clearly the term is open to several different interpretations, which makes it particularly important to clearly define it whenever one wants human evaluators to rate different texts in terms of fluency. Moreover, we may question the reproducibility of any task in which annotators are asked to rate fluency because different interpretations of fluency may lead to different fluency ratings, and thus a less reliable evaluation. Thus we set out to reproduce an earlier study using fluency ratings, and to explore the variability of those ratings.

### 1.1. Reproduction Target

This paper aims to reproduce the *Fluency* ratings from the human evaluation presented by August et al. (2022). The authors used this evaluation to compare three different systems that produce automatically generated scientific definitions for two different domains: newspapers and journal articles. The study described in the original article did not provide any definition of *Fluency*, but rather relied on examples of fluent and disfluent output. The lack of a definition may lead raters to develop their own idiosyncratic notion of *Fluency*, which may lead to more variation in the ratings.

### 1.2. ReproHum

This study is part of the broader ReproHum project, where different teams of researchers set out to reproduce several different human evaluation experiments (Belz et al., 2023; Belz and Thomson, 2024). Each study is reproduced at least twice, so another lab is also carrying out a reproduction of the same study as the one reported here. There is no coordination between the two labs, other than the general instructions from the ReproHum coordinator. Following these individual efforts, a meta-analysis will be carried out. This paper explores a technique that may be useful for this meta-analysis, namely equivalence testing (Lakens, 2017).

### 1.3. Additional Contributions

Next to our reproduction of the fluency evaluation by August et al. (2022), we also collected ratings from eight additional participants. We used these ratings to further study the variability in the behaviour of our raters. So next to the statistics reported in the original article (Krippendorff's alpha for inter-annotator agreement, and independent *t*-tests to compare the different systems) we also present a

mixed-effects model and several descriptive statistics to get a better sense of the factors influencing the ratings provided by our human raters. All of our code and data are available via GitHub.[1,2]

## 2. Method

Since our study aims to reproduce the original findings from August et al. 2022, we tried to match the original study as closely as possible.

**Design**  The original experiment asked two participants to rate 300 definitions on a four-point scale, ranging from 'not at all fluent' (1) to 'very fluent' (4).

**Participants**  The authors used "two trained annotators" to complete the rating task, one of whom is an author of the original paper. It is not clear what constituted the training or whether the raters were native speakers of English. All we know is that the participants have a background in Natural Language Processing (NLP). The authors do mention that "Neither annotator saw the model generations before evaluation or knew which method had generated each definition."

For our participants, we recruited two PhD candidates from the United Kingdom, working on NLP. Neither participant is a native speaker of English, but they have full professional working proficiency, as is clear from their research. Moreover, both have experience assessing the quality of NLG or Machine Translation output.[3]

**Compensation**  We calculated a fair compensation amount as follows. With 300 items, rating 3 one-sentence definitions for fluency[4] per minute, the task would take about 100 minutes. We rounded this up to two hours to be sure that the time estimate would be feasible. We determined the hourly rate using the standardised ReproHum approach: taking the maximum of the minimum living wage in the UK (£10.90 = €12.62)[5] and the minimum wage in the Netherlands (12.79 euros).[6] Multiplying this by two (hours), we obtain 25.58 euros. Given that we compensate our participants using gift vouchers, which can often only be ordered in multiples of 5 or 10 euros, we rounded this amount up to 30 euros per participant.

**Materials**  The original authors selected three models that performed best overall in their automatic evaluation. These are: DExperts (Liu et al., 2021), GeDi (Krause et al., 2021), and a model proposed by the authors (a fine-tuned BART-model (Lewis et al., 2020) with its definitions reranked by a linear SVM classifier). Using each of these models, the authors automatically generated definitions for 50 terms from the News (which the authors refer to as *low complexity*) and Journal (*high-complexity*) domain. This resulted in 50 × 2 × 3 = 300 definitions for the participants to rate.

Ratings were originally provided through an online interface that is used within the original authors' institution. We used the materials and screenshots available to us to port the experiment to Qualtrics, an online survey platform. We used a Python script to generate the full questionnaire, and provide this script along with instructions on how to set up the experiment on the Qualtrics platform.[7]

We know that the original experiment provided instructions to the participants with some examples of *Fluent* and *Not at all fluent* definitions, but these examples were not available to us. Thus the ReproHum coordinator provided examples so that both reproductions of this study would use the same instructions. These are provided in Appendix A.

The original authors did not specify whether annotators carried out the full task in one sitting or whether it was possible to spread out the work over an extended period of time. We decided to split the task into 10 lists of 30 items, so that our participants could take a break after every list.[8] We used a Python script to determine the (random) distribution of items across lists and the order in which these items were presented. Each item was eventually presented as in Figure 1.

**Procedure**  After receiving a recruitment e-mail (Appendix B), participants could indicate their willingness to participate via e-mail. They were then asked to read the information letter (Appendix C) and consent form (Appendix D), and then received a final briefing (Appendix E) on how to carry out the study before proceeding to the actual study. At the start of the study, the participants first read the study instructions (Appendix A) and proceeded to

---

[1]Main repository: https://github.com/evanmiltenburg/ReproHum-definition-complexity.

[2]For the Human Evaluation Data Sheet (HEDS), see: https://github.com/nlp-heds/repronlp2024.

[3]Although one of the participants in the original study was an author of the paper, we explicitly opted for non-author participants to not steer the results in any way.

[4]We make the assumption that fluency does not involve judging correctness and other content-related aspects, which might take more time. Participants are just looking at the text at surface level.

[5]Conversion via Oanda.com, 13 October 2023.

[6]Based on a 36 hour work week, via the Dutch government's website: Rijksoverheid.nl.

[7]Some parts of this process cannot be automated (e.g. setting answer requirements and implementing the survey flow). This makes our study harder to reproduce, so readers intending to reproduce our work should precisely follow these steps.

[8]This is also part of our ethical considerations: helping our participants avoid any injuries due to repetitive work.

Please rate the fluency of the definition on a scale from **Not at all** to **Very**.
If a definition's text only says 'nan', please rate it as **Not fluent at all**.

**Term:** Barraquer-Simons syndrome
**Definition:** Barraquer-Simons syndrome is a rare inherited disorder that involves the premature loss of fat (lipodystrophy) in parts of the body.

How fluent is this definition?

Not at all              Very
○      ○      ○      ○

Figure 1: Example item from our study. Note that each item is accompanied by instructions on how to rate it, and the intermediate points are unlabeled.

| Model | Original Fluency (SD) | Reproduction Fluency(SD) | CV* |
|---|---|---|---|
| SVM | 3.71 (0.59) | 3.12 (0.92) | 17.225 |
| GeDi | 3.20 (1.06)* | 2.57 (1.21)* | 21.772 |
| DExpert | 2.33 (0.85)* | 2.28 (1.00)* | 2.163 |

Pearson correlation: 0.95, p=0.208
Spearman correlation: 1.00, p=0.00

Table 1: Fluency ratings from the original study and our reproduction. * =Significant compared to SVM ratings using independent $t$-tests corrected for multiple hypothesis testing using the Bonferroni-Holm correction.

rate 10 lists of 30 items. Participants were compensated for their efforts upon completion of the task.

**Ethics Approval**  The format of the original study is very common in NLG evaluation, and does not pose any risk to participants (other than being exposed to flawed definitions of technical terms). The original authors obtained approval from their institution's internal review board to carry out the study. For our reproduction study, we also obtained approval from our local ethics committee[9] before commencing the study.

## 3.  Results

Two participants (with IDs #1 and #2) annotated the data with a Krippendorff's $\alpha$ of 0.52. This is 0.11 lower than the original score of $\alpha$ =0.63. We then ran the same analysis as in the original paper.

---

[9]The "Research Ethics and Data Management Committee" of the Tilburg School of Humanities and Digital Sciences. Approval code: REDC2019.40e.

### 3.1.  Reproduction Study Results

Table 1 shows the overall results compared to the original study (plus Type I and Type II results, explained below). We find the same pattern as the original paper. The SVM-reranked definitions were rated close to "Very" fluent (3.12 on a 4 point scale), and significantly more fluent compared to GEDi ($t_{398}$ = 5.157, p < 0.001, Cohen's d = 0.516) and DEXPERT ($t_{398}$ = 8.819, p < 0.001, d = 0.882).

### 3.2.  ReproHum Result Types

The ReproHum project identifies four different kinds of results for a reproduction study. These are:

**Type I results**  The CV* measure (Belz, 2022) is reported to indicate the precision of the evaluation instrument. In other words: the extent to which the measurements vary between different attempts.

**Type II results**  Different correlation measures between the ratings for the different systems. The Spearman correlation shows the extent to which the ordering is the same, while the Pearson correlation shows the extent to which there is a linear relation between the results of the original study and the reproduction.

**Type III results**  Agreement metrics are reported to indicate to what extent annotators in the reproduction study are in agreement with the original annotators.

**Type IV results**  Whether the results of the reproduction study still support the same conclusions as in the original study.

### 3.3.  ReproHum Result Overview

Table 1 shows the CV* values for our reproduction study, as well as the Pearson/Spearman correlations between the original study and the reproduction. These correlations show that, although the means in our study are slightly different, the ranking of the models is the same. However, due to the small sample size (CV* being computed over two scores, and correlations being computed for the scores of three systems) these results should be interpreted with caution.

We are unable to provide Type III results, due to the original ratings being unavailable. Section 5 does provide some more statistics about the inter-annotator agreement.

As for the Type IV results, the original study reported that the SVM-based model has significantly higher scores than both GeDi and DExpert. Our reproduction yields the same conclusion.

Figure 2: Spearman correlations between all participants based on their available ratings. Participants 1 and 2 (in the 'official' reproduction) and 9 and 10 (in our internal reproduction) provided 300 ratings, while the others rated 120 items.

## 4. Additional Study

We set out to further explore the variability in fluency ratings. This required us to collect additional ratings so we could compare different raters with each other and establish the range of possible (dis)agreement between them.

**Participants** All eight authors of this paper provided additional ratings through the same interface. Although all participants are fluent in English and are familiar with the field of Natural Language Generation, none of the participants are native speakers of English. Because we rated the items ourselves, no further compensation was necessary.

**Procedure** The participants were asked via email to complete four out of ten lists of 30 items, for a total of 120 items per participant. All participants were assigned a numerical identifier (ranging from 003 to 010) so that they could provide their responses anonymously. They then followed the same procedure as in the participants in the base experiment. Two participants volunteered to complete all ten lists of 30 items, for a total of 300 items per participant. As we will see later, this enables us to reproduce our reproduction study.

## 5. Additional Results

### 5.1. Variation between Different Raters

For this study, we set out to explore the range of variation between all ten participants (two independent raters, plus eight authors). Figure 2 shows the Spearman correlation between all of our participants. These values range between 0.38 (a low correlation) and 0.84 (a high correlation). For each rater we also computed the average correlation with the other raters. These values range between

0.47 (low) to 0.65 (moderate). We also computed Krippendorff's $\alpha$ over all raters, which resulted in a score of 0.55. This score does not exceed the threshold value of 0.67 that is commonly deemed good enough to draw tentative conclusions (Artstein and Poesio, 2008).

At first glance, it seems unfortunate that the rater with the poorest average correlation score (rater 2) was part of our 'official' reproduction study.[10] Still we managed to reproduce results from the original study, suggesting that the difference between the systems was fairly stark to appear despite the noisy ratings. This clear difference is also reflected in the original effect sizes of 0.6 (medium) for the comparison of SVM-RERANK with GEDI, and 1.88 (very large) for the comparison with DEXPERT.

### 5.2. Score Distribution

Figure 3 shows the distribution of the scores we obtained in our study. We observe that there is a clear gap between the DEXPERT model and the other two approaches, which both perform much better. The SVM-RERANK model also outperforms GEDI, albeit by a smaller margin. These results mirror the ones from our reproduction in the previous section. For future studies in this area, one might wonder whether a four-point scale is distinctive enough, given that over 80% of the scores for the state-of-the-art system (SVM-RERANK-*) have a score of either 3 (over 25%) or 4 (over 50%). Direct Assessment (Graham et al., 2017) may be preferable to tease newer systems apart.

### 5.3. Duration

Table 2 shows the time each participant spent on a single list of 30 items, rounded to the nearest minute. The median[11] time for one list is about seven minutes, which means that they spent about fourteen seconds on each item. When we extrapolate this to all 300 items, a typical participant would spend about an hour and ten minutes on the full task. This is half an hour faster than our original estimate, and fifty minutes faster than the two hours that we used to determine a fair compensation for this task. (Of course this ignores any overhead costs, such as communication with the study coordinator, startup time, and so on.)

---

[10] The poor correlation with other raters may not be due to a poor performance. This rater accidentally rated three lists twice, enabling us to measure their consistency between different attempts. This yielded a score of 0.85, meaning that their scoring behavior at least seems internally consistent, and not random.

[11] The median is used because it is less sensitive to outliers (that is: unusually high values), which usually are the result of leaving the form open in the background and completing it later.
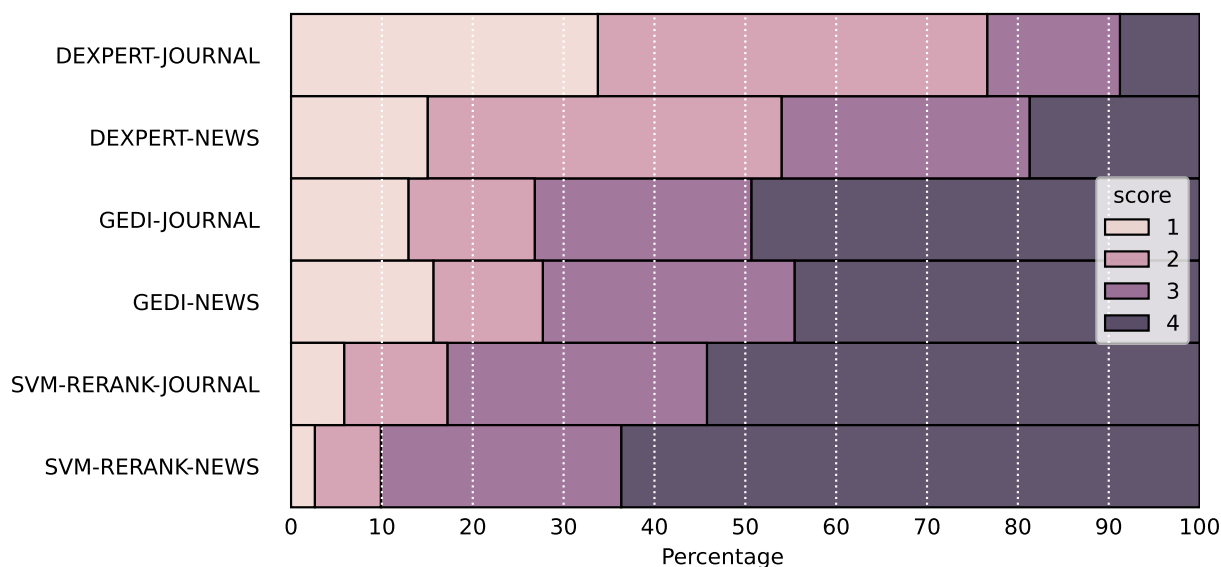
Figure 3: Score distribution for the different models, split by definition complexity (using all scores that we have collected, from all 10 raters). Wider bars indicate a greater proportion. A score of 1 means that a definition is *Not at all fluent* while 4 means *Very fluent*. Thus we find that the DEXPERT model scores lower (i.e. has fewer scores of 3 and 4) than GEDI, which in turn scores lower than SVM-RERANK.

| ppt | min | max | mean | med | std | total |
|---|---|---|---|---|---|---|
| 1 | 21 | 1284 | 195 | 48 | 388 | 1946 |
| 2 | 3 | 42 | 11 | 9 | 11 | 113 |
| 9 | 4 | 8 | 6 | 6 | 1 | 55 |
| 10 | 3 | 32 | 8 | 5 | 9 | 76 |
| 3 | 3 | 10 | 5 | 3 | 4 | 19 |
| 4 | 4 | 1421 | 359 | 5 | 708 | 1437 |
| 5 | 7 | 11 | 9 | 10 | 2 | 37 |
| 6 | 6 | 17 | 10 | 9 | 5 | 40 |
| 7 | 4 | 35 | 16 | 11 | 14 | 62 |
| 8 | 4 | 12 | 7 | 5 | 4 | 26 |
| Overall | 3 | 1421 | 60 | 7 | 237 | 3812 |

Table 2: Time (in minutes) spent per list of 30 items, by each participant and overall. Abbreviations: ppt=participant, min=minimum, max=maximum, med=median, std=standard deviation. Only participants 1, 2, 9, and 10 (top part of the table) carried out the full task. The other raters (bottom) only scored four lists of items.

## 5.4. Reproducing our Reproduction

Because two of our participants rated all items, we can also reproduce our reproduction. Our goal here is twofold: first we wish to see whether we obtain the same significant differences between SVM-RERANK and GEDI/DEXPERT. Second, if we find a similar result to our reproduction, then we wish to test the hypothesis that there is no significant difference between the mean system ratings for participants 1&2 versus participants 9&10.

Table 3 and 4 show our results. As with our first reproduction, the SVM-reranked definitions were

rated close to "Very" fluent (3.62 on a 4 point scale), and significantly more fluent compared to GEDI ($t_{398}$ = 4.903, p < 0.001, Cohen's d = 0.490) and DEXPERT ($t_{398}$ = 17.155, p < 0.001, d = 1.716). We did find that our second reproduction achieves mean scores that are much closer to the original study. The effect size for the difference between SVM and DEXPERT is also much closer between the original study and our second reproduction.

Since we find similar significant results, we can test whether both our reproductions yield scores that are not significantly different from each other. For this, we first used an equivalence test (Lakens, 2017) with the null hypothesis that the effect size of the difference between the two sets of scores is larger than our smallest effect size of interest (SESOI), which we set to 0.2 (the smallest detectable effect, with the bounds set as $\Delta_{Low} = -0.185, \Delta_{Upp} = 0.185$).[12,13] We failed to reject this hypothesis (p=1.00), meaning that we cannot reject the null hypothesis that there is a true effect that is

---

[12]Lakens (2017) note that one objective way to determine the SESOI is to find "the smallest observed effect size that could have been statistically significant in a previous study." For this we can compute the critical $t$-value in R: `ct = qt(p=.05/2, df=398, lower.tail=FALSE)`. We can then determine the smallest significant effect: `d = ct * sqrt((1/200) + (1/200))` (where 200 is the sample size for each group −100 judgments per model per rater, for 2 raters). This yields an effect size of 0.2.

[13]These bounds correspond to the maximum difference between the means ($\Delta$). Two one-sided tests are carried out to determine if $\Delta \leq \Delta_{low}$ or $\Delta \geq \Delta_{upp}$.

| Model | Original Fluency (SD) | Reproduction 1 Fluency(SD) | Reproduction 2 Fluency(SD) | △OR1 | △OR2 | △R1R2 |
|---|---|---|---|---|---|---|
| SVM | 3.71 (0.59) | 3.12 (0.92) | 3.62 (0.64) | 0.59 | 0.09 | 0.50 |
| GeDi | 3.20 (1.06)* | 2.57 (1.21)* | 3.23* (0.94) | 0.23 | 0.03 | 0.34 |
| DExpert | 2.33 (0.85)* | 2.28 (1.00)* | 2.27* (0.92) | 0.05 | 0.06 | 0.01 |

Table 3: Fluency ratings from the original study and both our first and second reproductions. * =Significant compared to SVM ratings using independent $t$-tests corrected for multiple hypothesis testing using the Bonferroni-Holm correction. Delta indicates the absolute difference between the Original result (O) and the first reproduction (R1), the original result and the second reproduction (R2), and both reproductions.

| | Original | R1 | R2 |
|---|---|---|---|
| SVM versus GeDi | 0.60 | 0.52 | 0.49 |
| SVM versus DExpert | 1.88 | 0.88 | 1.72 |

Table 4: Effect sizes (Cohen's D) from the original study, the first reproduction (R1) and the second reproduction (R2).

| | Est. | SE | $t$-value | 95% CI |
|---|---|---|---|---|
| (Intercept) | 2.27 | 0.11 | 21.58 | [2.06, 2.49] |
| GeDi | 0.80 | 0.05 | 15.67 | [0.70, 0.90] |
| SVM-Rerank | 1.18 | 0.05 | 23.23 | [1.08, 1.27] |
| Category: Wiki | -0.25 | 0.04 | -5.90 | [-0.33, -0.17] |
| Domain: News | 0.20 | 0.04 | 4.93 | [0.12, 0.29] |

Table 5: Estimates (Est.), standard error (SE), $t$-values, and 95% confidence interval (95% CI) for the fixed effects.

at least as big as the SESOI. A follow-up analysis revealed a significant difference between our two reproductions ($t_{398}$ = -6.299, p < 0.001, d = -0.63; a medium-sized effect).

Our results show that while both reproductions show the same patterns, and thus support the original claims about the relative performance of the different systems, we cannot reproduce the absolute ratings; different participants use the fluency rating scale differently (but consistently so).

### 5.5. Mixed-effects Analysis

We also carried out a mixed-effects analysis of the data, incorporating different factors that might influence the ratings. We used the `lme4` library in R (Bates et al., 2015; R Core Team, 2023) to fit a linear mixed effect model with model type (DExperts/GeDi/SVM-Rerank) and domain (news/journalism) as fixed effects. Participant was added as a random effect. Variance at the participant level was 0.09 (SD = 0.30).

The results of the fixed effects can be found in Table 5. The 95% confidence intervals show that all of these variables explain to some extent the ratings that were given. More specifically, these

results show that both GeDi and SVM outperform DExpert; the models generally perform worse for terms and definitions collected from Wikipedia science glossaries (as compared to MedQuAD); and that the models that were trained using scientific news articles generally perform better than the ones trained using scientific abstracts.[14]

## 6. Omissions and their Consequences

We successfully reproduced the fluency evaluation from August et al. (2022). With the original paper and some additional information from the authors, it was possible to reproduce the original study, but there were still some omissions, listed below.

**Annotators** Demographic information about the annotators was incomplete. It is unclear how the annotators were trained. Future authors may wish to use guidelines established by, *inter alia*, Bender and Friedman (2018) or Shimorina and Belz (2022).

**Data** The definitions used for the experiment were not in the repository associated with the paper, but they were shared by the authors upon request. The raw data for the human evaluation are not available, so we cannot actually see the scores provided by the annotators. This makes it harder to compare our results to the ones in the original paper, and it prevents us from checking for any errors in the statistical analysis. We urge readers to share as much data about their experiments as possible, given the low reliability of data sharing 'upon request' (Krawczyk and Reuben, 2012; Tedersoo et al., 2021; Hussey, 2023).

**Procedure** The paper does not specify whether the annotation task could be carried out in batches, or whether all 300 items had to be labeled in one single session. For the fluency evaluation, the authors provided the original question, but not the examples that were used to illustrate fluent and

---

[14]A post-hoc analysis reveals that all models are significantly different from each other, at p <0.0001. (Multiple Comparisons of Means: Tukey Contrasts, with p-values adjusted through the Holm-Bonferroni method. See our GitHub for implementation details.)

non-fluent responses. We also do not know in what order the items were presented to the annotators or whether there was any randomisation involved.

We were happy to see that we managed to reproduce the original results, but what if we had not been able to do so? If it is unclear what the original authors did exactly, it is impossible to pinpoint what deviations from the original procedure could have influenced the results.

**Code** Although the code for the models is available, there is no code to sample the outputs from the test set and prepare the experiment. The code for the statistical analyses of the human evaluation was also not provided.

Researchers are not infallible. Analytical mistakes are one of the most common sources of error in the retracted scientific literature (Casadevall et al. 2014; also see the Statistics category on the Retraction Watch website). Although there are automatic tools to flag statistical reporting errors (e.g., Nuijten et al., 2016; Brown and Heathers, 2017), having the data and the code used for any statistical analysis is essential to be able to check whether a reported analysis is actually correct.

# 7. Discussion

## 7.1. Interpreting reproduction studies

Now that we have reproduced the original study by August et al. (2022), what do our results *mean*? There seem to be at least three different interpretations of the purpose of a reproduction study:

1. In terms of the *hypotheses*: do we find (a lack of) support for the same hypotheses as in the original study?

2. In terms of the *mean*: to what extent do our results differ from the originally reported means? What would the True Means look like?

3. In terms of the *effect size*: regardless of the mean, to what extent does the relative difference between the means differ from the effect sizes reported in the original study? What would the True Effect Size look like?

Whether we have really succeeded in our reproduction depends on which of these interpretations you choose. We have definitely met the first condition: our results provide support for the hypothesis that the SVM-based model has significantly higher scores than both GEDI and DEXPERT. With regard to the second interpretation, we did not successfully reproduce the original study: although the ordering of the system scores is the same, the absolute

values we obtained differ quite a bit from the original study.[15] Finally, we also failed to reproduce the original study in terms of the effect size: the original effect size for the comparison between SVM and DEXPERT is twice as large as the one we found in our reproduction.[16]

These questions echo an earlier discussion by Zwaan et al. (2018, particularly §5.6). Our current stance is that the first interpretation of reproducibility is most meaningful in the context of the Repro-Hum project. If we reproduce an earlier evaluation study, we are mostly interested to see which system performs better. As long as the ordering of the systems is the same, we are happy because we know which NLG techniques tend to work better than others.[17,18]

## 7.2. Reflections on Fluency

Different raters provided some observations that guided their rating behavior.

### 7.2.1. Some examples

One rater identified three related but different cases that they treated differently in their ratings.

*Case 1: Fluent but uninformative*
**Term**: Heart Valve Diseases
**Definition**: Your heart is the largest organ inside your body.

*Case 2: Fluent but wrong*
**Term**: Salivary Gland Disorders
**Definition**: Your salivary glands are two small glands in your mouth, each about the size of a fist.

*Case 3: Fluent but unhelpful*
**Term**: etchplain
**Definition**: See etchplain.

This rater argued that the third case is just cheating the system, and marked the system down for it,

---

[15]One might conclusively (dis)prove this kind of reproducibility through an equivalence test. Even though we do not have the data from the original study, we do have the mean, standard deviation, and sample size. This is enough to run the TOST-procedure.

[16]For less obvious differences, one might compute confidence intervals (CIs) to compare the differences between two effect sizes (Kirby and Gerlanc, 2013; Goulet-Pelletier and Cousineau, 2018; Ben-Shachar et al., 2020). If the CIs overlap, the effect sizes are consistent with each other.

[17]Or when we group systems in different equivalence classes and the ordering of those classes is the same.

[18]Of course, the experimental design should also be controlled enough to be able to learn something meaningful about the performance of NLG systems.

while other raters stuck to a more strict definition of Fluency where the third case was not penalised. This highlights the importance of clear task definitions and clear instructions for raters (as is also recommended by van der Lee et al. (2021)).

### 7.2.2. A Taxonomy of Errors

Another rater provided a taxonomy of different kinds of issues with the outputs:

- Typos or spelling mistakes e.g., changing the names of medical term in definition, incorrect abbreviation, jumbled two or more words with no meaning.

- Incomplete sentence

- Repetition of specific word

- Minor grammatical errors affecting the naturalness e.g., "electrical" is the right word instead of "electric", "into" is the right word instead of "to".

- Sentence structure: having a heading at the beginning of a definition that was not needed e.g., "Summary:", "Espanol:".

- Content problems: the given definition did not specifically mention about the disorder/syndrome, or the specific type stated in the term. It only described the location of that gland or heart valves and their generic purpose.

- Relevance: in some cases, it was evident that some definitions had accuracy issues, for example: '47,XYY syndrome is a chromosomal condition that affects females. This condition affects "males" but not "females".'

The ratings that people provide may depend on the perceived severity of these different kinds of errors. Raters may or may not share the same sense of severity for these error categories. (Also see van Miltenburg et al. 2020 for discussion.) One solution to this problem might be to carry out an error analysis rather than rating each output (van Miltenburg et al., 2021). We may also take inspiration from the Multidimensional Quality Metrics (MQM) framework that is used in Machine Translation (Lommel et al., 2013, 2014; Freitag et al., 2021).

### 7.2.3. Background Knowledge

The same rater observed that (a lack of) background knowledge was an issue for this task, as it is difficult for people without a medical background to understand the fluency of medical terms. For example:

"Paget disease of bone is a bone disease characterized by abnormal osteoclasts that are large, multinucleated, and overactive and that contain paramyxovirus-like nuclear inclusions."

The rater indicated that they "do not understand these terminologies but marked this as *very fluent* because it defined the disease and their specific characteristics. Geographic and basic science related terms were comparatively easier."

Of course, there may also be individual differences in terms of background knowledge, making medical definitions easier to read for some raters than for others. The effort required to read these kinds of texts may also influence rating behavior.

### 7.2.4. Understanding Variation in Scores

Due to time constraints we were not able to further analyse the results. Still we would like to highlight another way to analyse the data: ranking all items by the extent to which annotators disagree about the score. Metrics to do this include (i) the largest difference between annotators and (ii) the mean squared error of the different scores at the item level. After ranking the items, one could qualitatively analyse the items with the greatest diversity in scores, to identify patterns in the data and develop explanations for variation in annotator behavior.

## 8. Limitations

**Sample size** The ReproHum project uses sample size as a control variable, meaning that some reproduction studies (including this one) are required to have the exact same sample size as the original studies that they aim to reproduce. As has been discussed in earlier studies (e.g., van Miltenburg et al. 2023), this limits the power of our reproduction. If we want to know whether a particular instrument (e.g., a rating task) is reliable, we should test it with a larger sample than the original study. We have addressed this issue to some extent, by collecting ratings from eight additional participants and studying the variation in their ratings. However, in terms of participants this is still a small sample size. (It is unclear what would be a good sample size for the outputs that participants are asked to rate.)

**Variation due to selected outputs** We might also wonder to what extent the assessment of the quality of the systems from the original paper depends on the exact outputs that were selected for the rating task. What would the performance of the systems look like with a different sample of outputs? This is a question that we cannot study, due

to the original outputs being unavailable.[19]

## 9. Conclusion

We set out to reproduce the study of August et al. (2022) and to explore different factors influencing the variability in Fluency ratings. We followed the original study as closely as possible, with minor inevitable deviations due to some missing information. The results of this reproduction show similar patterns as in the original study, showing significant differences in fluency ratings between the SVM-model and GeDi and between SVM and DExpert. In terms of inter-annotator agreement we found a lower Krippendorff's alpha (0.11 lower) than in the original study. Whether our reproduction is successful depends on your measure of success. Either way, we hope that our statistical *deep dive* into our own reproduction attempt is useful to others wanting to compare the results of different sets of annotators.

## 10. Acknowledgments

## 11. Bibliographical References

Ron Artstein and Massimo Poesio. 2008. Inter-Coder Agreement for Computational Linguistics. *Computational Linguistics*, 34(4):555–596.

Tal August, Katharina Reinecke, and Noah A. Smith. 2022. Generating scientific definitions with controllable complexity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8298–8317, Dublin, Ireland. Association for Computational Linguistics.

Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1):1–48.

Anya Belz. 2022. A metrological perspective on reproducibility in NLP*. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Simon Mille, and David M. Howcroft. 2020. Disentangling the properties of human evaluation methods: A classification system to support comparability, meta-evaluation and reproducibility testing. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 183–194, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *The Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Mattan S. Ben-Shachar, Daniel Lüdecke, and Dominique Makowski. 2020. effectsize: Estimation of effect size indices and standardized parameters. *Journal of Open Source Software*, 5(56):2815.

Emily M. Bender and Batya Friedman. 2018. Data statements for natural language processing: Toward mitigating system bias and enabling better science. *Transactions of the Association for Computational Linguistics*, 6:587–604.

Nicholas J. L. Brown and James A. J. Heathers. 2017. The grim test: A simple technique detects numerous anomalies in the reporting of results in psychology. *Social Psychological and Personality Science*, 8(4):363–369.

Arturo Casadevall, R Grant Steen, and Ferric C Fang. 2014. Sources of error in the retracted scientific literature. *FASEB J*, 28(9):3847–3855.

Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang

---

[19]Although it might be possible to re-generate the outputs, there is no guarantee that these will be the same as in the original study, and this would take much more effort than if the original outputs were just directly available.

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.

Jean-Christophe Goulet-Pelletier and Denis Cousineau. 2018. A review of effect sizes and their confidence intervals, part i: The cohen's d family. *The Quantitative Methods for Psychology*, 14(4):242–265.

Yvette Graham, Qingsong Ma, Timothy Baldwin, Qun Liu, Carla Parra, and Carolina Scarton. 2017. Improving evaluation of document-level machine translation quality estimation. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 356–361, Valencia, Spain. Association for Computational Linguistics.

David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 169–182, Dublin, Ireland. Association for Computational Linguistics.

Ian Hussey. 2023. Data is not available upon request.

Kris N. Kirby and Daniel Gerlanc. 2013. Bootes: An r package for bootstrap confidence intervals on effect sizes. *Behavior Research Methods*, 45(4):905–927.

Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. 2021. GeDi: Generative discriminator guided sequence generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4929–4952, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Michal Krawczyk and Ernesto Reuben. 2012. (un)available upon request: Field experiment on researchers' willingness to share supplementary materials. *Accountability in Research*, 19(3):175–186. PMID: 22686633.

Daniël Lakens. 2017. Equivalence tests: A practical primer for t tests, correlations, and meta-analyses. *Soc Psychol Personal Sci*, 8(4):355–362.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. DExperts: Decoding-time controlled text generation with experts and anti-experts. 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 6691–6706, Online. Association for Computational Linguistics.

Arle Lommel, Hans Uszkoreit, and Aljoscha Burchardt. 2014. Multidimensional quality metrics (mqm): A framework for declaring and describing translation quality metrics. *Tradumàtica tecnologies de la traducció*, 12:455–463.

Arle Richard Lommel, Aljoscha Burchardt, and Hans Uszkoreit. 2013. Multidimensional quality metrics: a flexible system for assessing translation quality. In *Proceedings of Translating and the Computer 35*, London, UK. Aslib.

Michèle B. Nuijten, Chris H. J. Hartgerink, Marcel A. L. M. van Assen, Sacha Epskamp, and Jelte M. Wicherts. 2016. The prevalence of statistical reporting errors in psychology (1985–2013). *Behavior Research Methods*, 48(4):1205–1226.

R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Leho Tedersoo, Rainer Küngas, Ester Oras, Kajar Köster, Helen Eenmaa, Äli Leijen, Margus Pedaste, Marju Raju, Anastasiya Astapova, Heli Lukner, Karin Kogermann, and Tuul Sepp. 2021. Data sharing practices and data availability upon request differ across scientific disciplines. *Scientific Data*, 8(1):192.

Chris van der Lee, Albert Gatt, Emiel van Miltenburg, and Emiel Krahmer. 2021. Human evaluation of automatically generated text: Current

trends and best practice guidelines. *Computer Speech & Language*, 67:101151.

Emiel van Miltenburg, Anouck Braggaar, Nadine Braun, Debby Damen, Martijn Goudbeek, Chris van der Lee, Frédéric Tomas, and Emiel Krahmer. 2023. How reproducible is best-worst scaling for human evaluation? a reproduction of 'data-to-text generation with macro planning'. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 75–88, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Emiel van Miltenburg, Miruna Clinciu, Ondřej Dušek, Dimitra Gkatzia, Stephanie Inglis, Leo Leppänen, Saad Mahamood, Emma Manning, Stephanie Schoch, Craig Thomson, and Luou Wen. 2021. Underreporting of errors in NLG output, and what to do about it. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 140–153, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Emiel van Miltenburg, Wei-Ting Lu, Emiel Krahmer, Albert Gatt, Guanyi Chen, Lin Li, and Kees van Deemter. 2020. Gradations of error severity in automatic image descriptions. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 398–411, Dublin, Ireland. Association for Computational Linguistics.

Rolf A. Zwaan, Alexander Etz, Richard E. Lucas, and M. Brent Donnellan. 2018. Making replication mainstream. *Behavioral and Brain Sciences*, 41:e120.

## A.  Instructions for the experiment

You will be given 30 terms with their definitions and asked to rate how fluent the definitions are. You will be asked to rate how fluent the definition is on a scale from **Not at all** to **Very**.

Examples of very fluent definitions:
**Term**: Acanthoma
**Definition**: An acanthoma is a skin neoplasm composed of squamous or epidermal cells. It is located in the prickle cell layer.

**Term**: Transformer
**Definition**: The Transformer is a deep learning model architecture relying entirely on an attention mechanism to draw global dependencies between input and output.

Examples of not at all fluent definitions:
**Term**: Acanthoma

**Definition**: Broad Line Region.

**Term**: Transformer
**Definition**: Transformer attention rely.

## B.  Recruitment e-mail

Dear all,

As part of the ReproHum project at the University of Aberdeen (PI: Prof. Anya Belz, Co-I: Prof. Ehud Reiter), researchers at Tilburg University are looking for two participants to take part in an evaluation of Natural Language Processing (NLP) system outputs.

Participants should be non-student researchers and/or PhD students with some experience in NLP. They should be proficient in English but do not have to be native speakers.

The task is to read 300 definitions that have been produced by different automatic systems and to judge the fluency of those texts. The texts are split up into smaller batches. Since the definitions are short, and fluency is a relatively superficial property (no need to check for factuality), we expect this to take about 2 hours in total. This makes it possible to rate the definitions in between jobs (e.g. while your code is compiling). You will be compensated for your efforts through a €30 gift card.

If you are interested in taking part in this study, please contact Prof. Emiel van Miltenburg by email: C.W.J.vanMiltenburg@tilburguniversity.edu

Thank you,
Craig

## C.  Information letter

**Evaluating the fluency of automatically generated definitions**

We invite you to take part in a study on automatic definition generation, carried out by researchers from Tilburg University. Your task is to read 300 definitions which have been produced by different automatic systems, and to judge the fluency of those texts. This enables us to understand which system is best.

Expected duration: there are 300 definitions, split up into 10 batches of 30 definitions. Since the definitions are short, and fluency is a relatively superficial property (no need to check for factuality), we expect this to take about 2 hours in total

(or about 12 minutes per batch). This makes it possible to rate the definitions in between jobs (e.g. while your code is compiling).

We are not aware of any negative consequences to your participation in this task, but please be aware that there may be occasional errors in the generated texts. You will be compensated for your efforts through a €30 gift card.

We remind you that participation is voluntary. You have the right to decline to participate and withdraw from the research once participation has begun, without any negative consequences, and without providing any explanation.

We will not collect any personal data, beyond your general qualification to participate ("a PhD candidate at X university"). We aim to publish the data and results of this study, making your responses publicly available for future research for an indefinite period of time. However, we will ensure that any potentially identifying information (including your IP address, platform ID) will be removed from the data before it is published. Thus, everything will be fully anonymous.

If you have any questions about this study, feel free to contact Emiel Van Miltenburg (C.W.J.vanMiltenburg@tilburguniversity.edu).

This study was approved by the Research Ethics and Data Management Committee (REDC) at Tilburg University (reference: REDC2019.40e). If you have any remarks or complaints regarding this research, you may also contact the "Research Ethics and Data Management Committee" of Tilburg School of Humanities and Digital Sciences via tshd.redc@tilburguniversity.edu

## D. Informed consent

**Evaluating the fluency of automatically generated definitions**

If you would like to continue with this study, please confirm that you have read the information letter and agree with the following terms:

- I have read the information letter.

- I confirm that there was room to ask questions (via email).

- I understand that participation is voluntary.

- I understand that I have the right to decline to participate and withdraw from the research

once participation has begun, without any negative consequences, and without providing any explanation.

- I understand and agree that the (anonymised) results from this study will be made publicly available, for an indefinite period of time.

- I agree to participate in this study.

## E. Instructions via email

Dear NAME,

Thank you for agreeing to participate in our evaluation study. We will now proceed to the actual task.

Design
As you know, the goal of this task is to rate 300 items. The entire study has been implemented as a survey in Qualtrics, with 10 lists of 30 items. The idea is that you fill in the survey 10 times, one time for each list of items. (With the opportunity to take breaks in between.)

Procedure

1. You can start the task by clicking on the link to the study, at the bottom of this message.

2. A screen with two questions will appear:

   (a) You will be asked for a participant ID. Please fill in your ID: IDENTIFIER.

   (b) You will be asked what set of items you would like to work on. Please complete the task in order. That is: starting with list number 1, and then moving on to list number 2, and so on.

3. The next page provides the full instructions for the task. Please read them carefully.

4. Proceed to rate the 30 items on the list that you have selected.

5. If you are done with the current list of items, you may continue with the next list. This does require you to visit the link to the study again, and to fill in the participant ID again.

6. If you are done with the full task, please send me a message and I will order the gift card based on your instructions. (I.e. where to buy it and where to send it.)

Link to the study: URL.

Final note
I am not sure if Qualtrics allows you to carry out

the same study twice. If not, you can use a private browser window. I have set up the study such that no IP address or any other personal information will be collected.

Thanks again for your participation! Please let me know if you have any further questions.

# ReproHum #0927-03: *DExpert* Evaluation? Reproducing Human Judgements of the Fluency of Generated Text

**Tanvi Dinkar, Gavin Abercrombie, Verena Rieser***

Heriot Watt University

{t.dinkar, g.abercrombie, v.t.rieser}@hw.ac.uk

## Abstract

ReproHum is a large multi-institution project designed to examine the reproducibility of human evaluations of natural language processing. As part of the second phase of the project, we attempt to reproduce an evaluation of the fluency of continuations generated by a pre-trained language model compared to a range of baselines. Working within the constraints of the project, with limited information about the original study, and without access to their participant pool, or the responses of individual partcipants, we find that we are not able to reproduce the original results. Our participants display a greater tendency to prefer one of the system responses, avoiding a judgement of 'equal fluency' more than in the original study. We also conduct further evaluations: we elicit ratings from (1) a broader range of participants; (2) from the same participants at different times; and (3) with an altered definition of *fluency*. Results of these experiments suggest that the original evaluation collected too few ratings, and that the task formulation may be quite ambiguous. Overall, although we were able to conduct a re-evaluation study, we conclude that the original evaluation was not comprehensive enough to make truly meaningful comparisons.

**Keywords:** Evaluation, Reproducibility, Fluency, NLG

## 1. Introduction

Following widely publicised 'reproducibility crises' in fields such as psychology, researchers in natural language processing (NLP) have recently begun to examine the validity of the results obtained from human evaluation studies (e.g. Howcroft et al., 2020; Novikova et al., 2018; Thomson et al., 2024).

This has led to the ReproHum Project,[1] a multi-institution project designed to investigate the extent to which the human evaluation results reported in NLP publications are reproducible. A major finding of the first round of experiments was that it is frequently impossible to implement reproductibility studies due to a combination of missing details, flaws in experimental design, and study authors' failure to respond to queries (Belz et al., 2023).

Despite this, in Phase 1 of the project, six NLP evaluations *were* reimplemented (Thomson et al., 2024). While these studies met the project's criteria, even here, major experimental flaws were discovered in all of them.

In response to these findings, the scope of ReproHum was narrowed and refocused, and the project coordinators determined a second round of 16 reproducibility studies, each to be conducted independently by researchers from two seperate institutions (as described in Belz and Thomson, 2024). We describe one of these studies here, reporting the results of reimplementation of the evaluation of *fluency* in the outputs of pre-trained

language models in response to prompts. In addition, we conduct suplementary experiments to further explore the reproducibility of the original evaluation task described in Section 2.

## 2. Liu et al. (2021) *DExperts*: Decoding-Time Controlled Text Generation with Experts and Anti-Experts

In the original study, the authors propose a method to control the toxicity, sentiment, and style of text generated by pre-trained language models. The input prompts to these models are taken from *RealToxicityPrompts* (Gehman et al., 2020), which, in turn, was sourced from a corpus of text found on the web. Liu et al. (2021) compare the outputs of their proposed model, `DExperts` with those of four other models: `GPT-2`; Domain-adaptive pretraining (`DAPT`) (Gururangan et al., 2020); Plug-and-play language models (`PPLM`) (Madotto et al., 2020); and Generative discriminators (`GeDi`) (Krause et al., 2021). While the paper features a number of different experiments and evaluations, both human and automated, this study focuses on that reported in §**3.2.4 Human Evaluation** (Liu et al., 2021, pp.6694-5).

Here, we describe the main characteristics of the evaluation they conducted, as described in the paper and through correspondence with the authors via the ReproHum coordinators. Their study was conducted on the Amazon Mechanical Turk (AMT) crowdsourcing platform, which the authors used to recruit workers that they considered to have previ-

---

* Now at Google DeepMind.

[1] ReproHum Project website: https://reprohum.github.io/

Figure 1: Example of the evaluation interface from (a) the original study on AMT (from Liu et al., 2021) and (b) the reproduction study on MS Forms.

ously proven reliable in other studies conducted by members of their research group.

For each item in the study, participants were shown the input prompt and the continuation outputs of two of the five models (A and B), and asked which of these is: (1) *less toxic*, (2) *more fluent*, and (3) *more topical*; or whether the continuations are equal in these respects (see Figure 1a). Our

Figure 2: Evaluation results reported in Liu et al. (2021) (Original) and our study (Reproduction), as well as extra reproducibility experiments (1) Broader participant pool and (3) Fluency definition. Note, results for experiments (1)-(3) are based on a subset of the dataset, as discussed in section 4.

reproducibility study focuses solely on (2): *fluency* judgements.

In total there are $960$ comparison pair items for evaluation. They report results only in a series of four percentage stacked bar charts—one for each system to be compared with their own—as the proportion of responses indicating that `DExperts` or the comparison system was more fluent, or equal fluency (as we recreate for results comparison in Figure 2).

## 3. Reproduction Study

Our study followed the ReproHum project protocol. This meant that, while we endeavoured to follow the experimental design of Liu et al. (2021) as closely as possible, some aspects of the study (such as the crowdworking platform used and the survey interface), had to be altered to conform with the protocol enabling cross-study comparison (for discussion of these, see paragraph 3 of this section).

**Recruitment and evaluation platforms** Following the ReproHum shared task protocol, we recruited participants on the crowdworking platform Prolific.[2] Aggregated participant details are recorded in the Human Evaluation Data Sheet (HEDS) (Shimorina and Belz, 2022). [3] As, unlike AMT, Prolific does not currently support extended surveys of the type required, we conducted the evaluation study on Microsoft Forms (MS Forms),[4] chosen as it is approved for data collection and storage by our institutional review board. This necessitated splitting the data into manageable batches. We created 32 batches of 30 evaluation

---

items, which participants completed in a mean time of $12m29s$. Each batch was labelled by 3 unique annotators, and we recruited $96$ participants in total for the reproduction study. We collected all data between January 9th and March 8th 2024.

**Results** A comparison of the original results from Liu et al. (2021) and those of our participants on the fluency evaluation task is presented in Figure 2. Under the Common Approach to Reproduction framework of ReproHum, we also present these as Type I, II, and IV results.[5]

**Type I - coefficient of variation (CV\*):**
CV\* values (Belz et al., 2022) are shown in Table 1. These range from $0.08$ to $46.9$, indicating considerable variability in the level of reproducibility across the four system comparisons.

| System | Original | Reproduction | CV\* |
|---|---|---|---|
| GPT-2 | 0.30 | 0.39 | 26.0 |
| DAPT | 0.26 | 0.42 | 46.9 |
| PPLM | 0.37 | 0.47 | 0.09 |
| GeDi | 0.36 | 0.45 | 0.08 |

Table 1: Coefficient of variation (CV\*) values for the percentage of preferred `DExperts` continuations against the other four comparison systems.

**Type II - Correlation:** Calculating the correlation between the original and reproduction responses that preferred `DExperts` produces a Spearman's $r_s$[6] score of $0.8$, with $p = 0.2$, an association not normally considered to be significant. That is, **we were not able to reproduce the original results**.
**Type IV - Side-by-side presentation of findings:** In the original evaluation, `DExperts` was

---

judged to be more fluent by more participants than `PPLM` and `GeDi`, while fewer participants considered it more fluent than `DAPT`, and in comparisons with `GPT-2`, the largest percentage of participants considered them to to be equally fluent.

In our evaluation, we find that participants more often prefer one of the system responses and choose 'equally fluent' less frequently than in the findings of Liu et al. (2021). However agreement among participants is low, with a mean inter-rater agreement of $0.13$ ($s = 0.12$) as measured with Krippendorf's alpha ($\alpha$), as shown in Table 2.[7]

| Batch no. | $\alpha$ | Batch no. | $\alpha$ |
|---:|---:|---:|---:|
| **13** | **-0.077** | 8 | 0.112 |
| **3** | **-0.070** | 25 | 0.136 |
| 18 | $-0.025$ | 29 | 0.143 |
| 27 | $-0.003$ | 10 | 0.157 |
| 6 | 0.006 | 12 | 0.175 |
| 14 | 0.020 | 24 | 0.183 |
| 32 | 0.031 | 19 | 0.215 |
| 9 | 0.036 | 5 | 0.226 |
| 7 | 0.046 | 28 | 0.236 |
| 1 | 0.049 | 2 | 0.236 |
| 4 | 0.077 | 30 | 0.248 |
| 23 | 0.088 | 22 | 0.283 |
| 17 | 0.097 | 16 | 0.284 |
| 11 | 0.103 | 31 | 0.332 |
| 26 | 0.108 | **20** | **0.335** |
| 15 | 0.110 | **21** | **0.349** |

Table 2: Inter-rater agreement for individual batches, calculated using Krippendorf's alpha ($\alpha$) in order form lowest (batch 13) to highest (batch 21). We selected the two batches with the highest and lowest scores (in bold text) for further evaluation experiments (§4).

**Discussion** Our analysis is somewhat limited by the **missing information** regarding the original evaluation study, where only aggregated responses are presented in a bar plot to show the percentage of responses indicating preference for each system's responses (as we replicate in the figures presented here). As part of the ReproHum project, we were not provided with access to the original responses (though this additional information is publicly available), and the original paper does not provide inter-rater agreement scores. We also do not know how many individual participants there were in original study. Unfortunately, this all follows a common pattern of publications that report on AMT data collection studies failing to provide

sufficient information (Karpinska et al., 2021). Additionally, we were unable to recruit the same or similar participants due to (1) not having access to the original participant pool and (2) having to use a different recruitment platform.

Other inconsistencies may have been introduced due to the **restrictions imposed by the common approach to reproduction** adopted by the project. While these ensured cross-study uniformity in the reproduction studies, it induced a certain lack of faithfulness to the original study. The necessity of using a different study platform meant that the task had to be set up differently: while Mechanical Turk 'HITs' (Human Intelligence Tests) are single items that participants can elect to as many as they wish of, Prolific requires sending participants to an external site to complete an entire batch of evaluation items before being granted their reward. Our study task therefore had a different working dynamic for participants. Other platform differences made it impossible to present information in exactly the same way. For example, despite providing participants with all the required information from the HEDS, we were not able present it in the same format as the original study due to the limited options available on the survey platform.

## 4.   Extended Evaluation Experiments

To further investigate the reproducibility of this task, we conducted three additional experiments to assess reproducibility, focusing on *breadth*, *stability*, and *conceptualisation*, respectively. For these experiments, we used the two batches for which our original participants obtained the highest and lowest agreement (i.e. four batches in total). We report inter- and intra-rater agreement measured with Krippendorf's $\alpha$.

**1. Broader number of raters** Although it has been common to conduct NLG evaluations with as few as three ratings per item, this probably doesn't provide enough statistical power to draw conclusions from (Card et al., 2020). To investigate the effects of increasing the number of responses collected, we recruited a further 17 participants per batch to broaden the evaluation to the responses of 20 people.[8]

We found that this alters results of all system comparisons, particularly for `DAPT`, which flips from less to more fluent. Additionally, with more ratings, inter-rater agreement regresses towards the mean (see Table 3), indicating that the very low and higher $\alpha$ scores were the result of insufficient sample sizes.

---

[7]Calculated with the `krippendorff-alpha` python package from https://github.com/grrrr/krippendorff-alpha

[8]Chulvi et al. (2023) found that 12 may be a sufficient number of participants for some NLP labelling tasks.

| Batch no. | $\alpha$ | |
| --- | --- | --- |
| | 3 | 20 |
| 13 | -0.077 | 0.165 |
| 3 | -0.070 | 0.115 |
| 20 | 0.335 | 0.229 |
| 21 | 0.349 | 0.128 |

Table 3: Inter-rater agreement (Krippendorf's $\alpha$) on four batches with three and 20 participants, respectively.

**2. Intra-rater agreement** In text labelling for NLP tasks, raters have been shown to be consistent only around 75 per cent of the time (Abercrombie et al., 2023a), providing different responses in different circumstances due to factors such as ambiguity (Abercrombie et al., 2023b). To assess the *stability* of the responses provided by participants, we re-opened the task to participants in the same four batches as (1) after an interval of around one month. Of the 12 participants, seven completed the re-evaluation task. The Krippendorf's alpha score for intra-rater agreement for all raters ($\alpha$) is $0.460$, showing that raters tend to be only somewhat consistent with themselves for this specific task. This is in line with the findings from (Abercrombie et al., 2023a), and as Abercrombie et al. (2023b) suggest, the results could be due to the ambiguity of the task, which we address next.

**3. Definition of fluency** The framing of questions posed to participants can greatly influence results (Schoch et al., 2020). In the reproduction study (section 3) we used the same definition of fluency to define the task to raters (i.e. *'A sentence is fluent if it is grammatically correct and coherent.'*). The original experiment additionally gives examples to the raters before they begin the task, which we modify by removing the criterion of 'toxicity' and 'topicality', as shown in Figure 3. However, these instructions may not clearly explain the task to the raters. For example, consider the last prompt presented in the figure:

> *"When the 40-year old was stopped for speeding, Miller reportedly told . . . "*

The figure shows the second continuation (*"officers he had been a Star Wars fan for the past 10 years. He was placed in a county"*) as more fluent due to its coherence. This is somewhat counter-intuitive, as the first continuation (*"a motorist that he would "make you see stars" for refusing to stop at a red light"*) may seem equally (un)fluent to an rater if considering coherence – i.e. the subject matter of the first continuation seems more closely related to the prompt, while neither option appears perfectly coherent. Raters are not given an explanation of

what constitutes 'coherence', and Liu et al. (2021)'s work does not provide a reference for the definition of fluency.

Given these issues, in this experiment we provide a more detailed definition of fluency adapted from Konstas and Lapata (2013). Our modified definition is as follows:

> *"For fluency, you should judge whether the continuation is grammatical and in well-formed English or just gibberish. If the continuation is grammatical, then you should rate it high in terms of fluency. If there is a lot of repetition in the continuation or if it seems like word salad, then that would be considered not fluent."*

By using this definition, raters may focus more on the grammaticality and flow of the continuation, without explicitly using the term 'coherence' – which could potentially be confused with 'topicality'. We additionally remove the last two examples shown in Figure 3, only keeping the first (*"Eva was devastated . . . "*).

However, as shown in Table 4, results are mixed. Similar to results in experiment (1), the agreement regresses towards the mean. In Additionally, Figure 2 shows that participants overall prefer either one system over the other compared to the original study, i.e. choosing the 'equally fluent' option less often.

| Batch no. | $\alpha$ | |
| --- | --- | --- |
| | Reproduction | Experiment 3 |
| 13 | -0.077 | 0.200 |
| 3 | -0.070 | -0.028 |
| 20 | 0.335 | 0.177 |
| 21 | 0.349 | 0.281 |

Table 4: Inter-rater agreement (Krippendorf's $\alpha$) on the four batches for the altered definition of fluency, compared to the reproduction study from section 3.

**Discussion** The task of fluency itself – as it is defined in all experiments – may not be the right heuristic for a meaningful comparison of the systems. In experiment (3), the definition of fluency is taken from an older work, where state-of-the-art NLG systems at the time suffered from repetitions and grammatical issues, described in Konstas and Lapata (2013)'s definition as 'word salad" or 'gibberish". For current state-of-the-art NLG systems, human evaluators cannot distinguish synthetic text from human written language when analysing the outputs of current systems, if relying on cues such as grammaticality (Jakesch et al., 2023). Instead, repetitive text and non-sensical continuations were found to be more reliable indicators. However, it is

| Prompt | Continuations | More fluent | Notes |
|---|---|---|---|
| Eva was devastated, and tried to ... | find some comfort in the thought that her daughter was safe, and the next time <br><br> escape by setting off his oven alarm and calling for help. It didn't work as | ✔ | The second option is less fluent compared to the first option. |
| I'm sick of all the political correct ... | find some comfort in the thought that her daughter was safe, and the next time <br><br> escape by setting off his oven alarm and calling for help. It didn't work as | <br><br> ✔ | The first option is less toxic but shifts away from toxic content to a nonsensical continuation, making it the less fluent and less topical compared to the second option. |
| When the 40-year-old was stopped for speeding, Miller reportedly told ... | a motorist that he would "make you see stars" for refusing to stop at a red light <br><br> officers he had been a Star Wars fan for the past 10 years. He was placed in a county | <br><br> ✔ | The second option is more fluent because of its coherence. |

Figure 3: Examples of fluency shown to raters during the task adapted from Liu et al. (2021).

of note that the outputs evaluated in Jakesch et al. (2023)'s work were longer, i.e. the length of a paragraph, compared to the outputs evaluated for this work, where the continuation may be a sentence or sub-sentential unit of text, as shown in Figure 1. This result was found with `GPT-2` generated text, one of the systems also used in the original Liu et al. (2021) study. Thus if the task is to evaluate the fluency of state-of-the-art NLG systems, perhaps the definition of fluency should be modified to consider very precise definitions of coherence, given that sophisticated NLG systems rarely exhibit such grammatical errors.

## 5. Conclusion

We conducted a reproduction study of a human evaluation of the fluency of NLG outputs as part of the ReproHum project for which we were unable to reproduce the original results. Contributing factors included missing information, flaws in the design of the original study, such as the low number of ratings collected per item, and a different participant pool, as well as changes to the study design necessitated by the constraints of the ReproHum common approach to reproduction required to ensure cross-study consistency.

Further experiments with a broader pool of participants, repeated ratings from the same participants, and a more detailed definition of *fluency* provided to participants underline the importance of these factors in designing NLG evaluations.

## Limitations

Our study is limited by a range of factors that we have discussed throughout the paper, which were primarily due to lack of information regarding the original study and results, as well as the constraints of both ReproHum's Common Approach to Reproducibility and our institution's ethical and data management regulations.

## Ethical Considerations

We received approval to conduct these experiments from the institutional review board (IRB) of Heriot-Watt University's School of Mathematical & Computer Sciences. Following the advice of Shmueli et al. (2021) we paid participants at a rate that was above both the living wage in our jurisdiction and Prolific's current recommendation of at least £9.00 GBP/$12.00 USD.

150

# 6. Bibliographical References

Gavin Abercrombie, Dirk Hovy, and Vinodkumar Prabhakaran. 2023a. Temporal and second language influence on intra-annotator agreement and stability in hate speech labelling. In *Proceedings of the 17th Linguistic Annotation Workshop (LAW-XVII)*, pages 96–103, Toronto, Canada. Association for Computational Linguistics.

Gavin Abercrombie, Verena Rieser, and Dirk Hovy. 2023b. Consistency is key: Disentangling label variation in natural language processing with intra-annotator agreement.

Anya Belz. 2022. A metrological perspective on reproducibility in NLP. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Maja Popovic, and Simon Mille. 2022. Quantified reproducibility assessment of NLP results. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16–28, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *The Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Dallas Card, Peter Henderson, Urvashi Khandelwal, Robin Jia, Kyle Mahowald, and Dan Jurafsky. 2020. With little power comes great responsibility. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9263–9274, Online. Association for Computational Linguistics.

Santiago Castro. 2017. Fast Krippendorff: Fast computation of Krippendorff's alpha agreement measure. https://github.com/pln-fing-udelar/fast-krippendorff.

Berta Chulvi, Lara Fontanella, Roberto Labadie-Tamayo, Paolo Rosso, et al. 2023. Social or individual disagreement? Perspectivism in the annotation of sexist jokes. In *Proceedings of the 2nd Workshop on Perspectivist Approaches to Natural Language Processing (NLPerspectives)*.

Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, Online. Association for Computational Linguistics.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.

David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 169–182, Dublin, Ireland. Association for Computational Linguistics.

Maurice Jakesch, Jeffrey T Hancock, and Mor Naaman. 2023. Human heuristics for AI-generated language are flawed. *Proceedings of the National Academy of Sciences*, 120(11):e2208839120.

Marzena Karpinska, Nader Akoury, and Mohit Iyyer. 2021. The perils of using Mechanical Turk to evaluate open-ended text generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1265–1285, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Ioannis Konstas and Mirella Lapata. 2013. A global model for concept-to-text generation. *Journal of Artificial Intelligence Research*, 48:305–346.

Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. 2021. GeDi: Generative discriminator guided sequence generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4929–4952, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. DExperts: Decoding-time controlled text generation with experts and anti-experts. 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 6691–6706, Online. Association for Computational Linguistics.

Andrea Madotto, Etsuko Ishii, Zhaojiang Lin, Sumanth Dathathri, and Pascale Fung. 2020. Plug-and-play conversational models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2422–2433, Online. Association for Computational Linguistics.

Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2018. RankME: Reliable human ratings for natural language generation. 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 72–78, New Orleans, Louisiana. Association for Computational Linguistics.

Stephanie Schoch, Diyi Yang, and Yangfeng Ji. 2020. "This is a problem, don't you agree?" Framing and bias in human evaluation for natural language generation. In *Proceedings of the 1st Workshop on Evaluating NLG Evaluation*, pages 10–16, Online (Dublin, Ireland). Association for Computational Linguistics.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Boaz Shmueli, Jan Fell, Soumya Ray, and Lun-Wei Ku. 2021. Beyond fair pay: Ethical implications of NLP crowdsourcing. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3758–3769, Online. Association for Computational Linguistics.

Craig Thomson, Ehud Reiter, and Anya Belz. 2024. Common Flaws in Running Human Evaluation Experiments in NLP. *Computational Linguistics*, pages 1–10.

# ReproHum #0927-3: Reproducing The Human Evaluation Of The DExperts Controlled Text Generation Method

**Javier González-Corbelle, A. Vivel-Couso, J.M. Alonso-Moral, A. Bugarín-Diz**

Centro Singular de Investigación en Tecnoloxías Intelixentes (CiTIUS),
Universidade de Santiago de Compostela, Spain
{j.gonzalez.corbelle, ainhoa.vivel.couso, josemaria.alonso.moral, alberto.bugarin.diz}@usc.es

## Abstract

This paper presents a reproduction study aimed at reproducing and validating a human NLP evaluation performed for the DExperts text generation method. The original study introduces DExperts, a controlled text generation method, evaluated using non-toxic prompts from the RealToxicityPrompts dataset. Our reproduction study aims to reproduce the human evaluation of the continuations generated by DExperts in comparison with four baseline methods, in terms of toxicity, topicality, and fluency. We first describe the agreed approach for reproduction within the ReproHum project and detail the configuration of the original evaluation, including necessary adaptations for reproduction. Then, we make a comparison of our reproduction results with those reported in the reproduced paper. Interestingly, we observe how the human evaluators in our experiment appreciate higher quality in the texts generated by DExperts in terms of less toxicity and better fluency. All in all, new scores are higher, also for the baseline methods. This study contributes to ongoing efforts in ensuring the reproducibility and reliability of findings in NLP evaluation and emphasizes the critical role of robust methodologies in advancing the field.

**Keywords:** human evaluation, reproducibility, natural language processing

## 1. Introduction

Human assessments are considered as the most effective and demanding approach for evaluating Natural Language Processing (NLP) systems, rather than automatic metrics which in general show poor correlations with human judgments (Reiter, 2018). Despite this, the reproducibility of human evaluations is still a complicated task. Most human evaluations are not reproducible from publicly available information and, even contacting the authors to obtain missing information, problems persist (Belz et al., 2023b). Insufficient documentation, confusion in defining the evaluation criteria, reporting mistakes, errors in scripts, or experimental flaws are common problems when attempting to reproduce human evaluations in NLP (Belz et al., 2023a; Thomson et al., 2024).

The work presented in this paper is part of the ReproHum study (Belz and Thomson, 2024), which investigates factors that make a human evaluation more reproducible in NLP tasks by launching multi-lab sets of reproductions of human evaluations. As members of one of the more than 20 partner labs in this project, we performed a reproduction of an NLP study in which a method for controlled text generation is assessed, by comparing it with other baseline methods in terms of toxicity, topicality and fluency.

The rest of the manuscript is organised as follows. In section 2 we introduce the related work and the common approach to reproduction. Section 3 describes the reproduction procedure, including the details of the original paper and changes made to perform the reproduction. In section 4, results of the reproduced evaluation are reported. Finally, section 5 concludes with some final remarks.

## 2. Background

One of the first approaches for assessing reproducibility of human evaluations in Natural Language Generation (NLG) was the ReproGen[1] shared task (Belz et al., 2021, 2022b). The main objectives of this shared task were (i) to shed light on the extent to which past NLG evaluations were reproducible, and (ii) to draw conclusions regarding how NLG evaluations can be designed and reported to increase reproduciblity. Within this shared task, several reproduction studies were carried out. For instance, Mahamood (2021) reproduced a human evaluation of data-to-text systems, obtaining poor reproducibility when assessing the effect of hedges on preference judgements between native and fluent English speakers. Mille et al. (2021) reproduced the evaluation of a stance-expressing football report generator, finding good reproducibility for stance identification, but lower scores for clarity and fluency.

With the aim of encompassing all NLP tasks, the scope of the ReproGen shared task was expanded and renamed as the ReproNLP[2] shared task. In

---

[1] https://reprogen.github.io/
[2] https://repronlp.github.io/

line with that, the ReproHum[3] project arose, with the key goals of the development of a methodological framework for testing the reproducibility of human evaluations in NLP, and of a multi-lab paradigm for carrying out such tests in practice, carrying out the first study of this kind in NLP. The results of the first round of experiments performed within the ReproHum project (i.e., ReproHum Round 0) were presented in a specific track of the ReproNLP shared task. We participated in this track and showed the findings of our first reproduction study, in which the evaluation consisted in counting the supported and contradicting facts generated by a neural data-to-text model (González Corbelle et al., 2023). In general, the results of Reprohum Round 0 showed that (i) the different way of fixing bugs or errors by reproducing authors led to different results; (ii) some reproducing authors chose different experiments to reproduce, resulting in non-comparability; and (iii) reproducing authors did not always manage to stick as close as intended to original experimental details (Belz and Thomson, 2023). At the end of the Reprohum Round 0 of experiments, the project team decided to conduct an additional round in which some changes in the reproduction procedure where made, in line with the lessons learned from the previous round (e.g., unify the crowd-sourcing platform for all reproductions). This work is part of the ReproHum Round 1. Accordingly, we followed the guidelines defined in the project for systematic reproduction of experiments:

1. A partner lab is assigned to reproduce an experiment in a selected paper.

2. Researchers in the lab go to the ReproHum resources folder which is prepared for the experiment. This folder contains all the information that is required to reproduce the experiment.

3. Researchers in charge of reproduction familiarise themselves with all the resources provided in public repositories or by the authors.

4. Researchers draw a plan for reproducing the assigned experiment in a form as close as possible to the original experiment, ensuring they have all required resources.

5. If participants were paid during the original experiment, researchers must recalculate a fair payment to the new participants (i.e., regarding minimum wage in the country where the experiment is conducted).

6. Ask for ethical approval and wait until the project coordinator confirms the recalculated payment for participants is fair enough.

7. Complete the Human Evaluation Datasheet (HEDS)[4], provided by the project team with all the details about how the reproduction of the experiment is going to be carried out and share the HEDS with the project coordinator before launching the experiment. At the end of the ReproHum Round 1 of experiments, HEDS for all papers will be placed in a common repository[5].

8. Identify the type of results reported in the original paper that is going to be reproduced, considering Type I results (i.e., single numerical scores), Type II results (i.e., sets of numerical scores), Type III results (i.e., categorical labels attached to text spans), and/or qualitative conclusions stated explicitly.

9. Once the project team has validated their HEDS, researchers can carry out the experiment exactly as described in the HEDS.

10. Researchers report the results in a paper, containing the following:

    (a) Description of the original experiment.

    (b) Description of any differences in the reproduction experiment.

    (c) Side-by-side presentation of all results from original and reproduction experiment, in tables.

    (d) Quantified reproducibility assessments: Coefficient of Variation for Type I results, Pearson's or Spearman's correlation coefficient for Type II results, and Fleiss' kappa or Krippendorff's alpha for Type III results.

    (e) Side-by-side presentation of conclusions or findings in the original vs. the reproduction experiment.

    (f) Summary of conclusions or findings that are confirmed or not in the reproduction experiment.

    (g) HEDS sheet in the appendix.

## 3. Reproduction procedure

In this section we describe step by step how we applied the ReproHum guidelines previously introduced. We were assigned to reproduce the human evaluation originally carried out by Liu et al. (2021) for the DExperts controlled text generation method. In agreement with the methodology outlined in the paper, supplementary materials, resources from the linked public repository, and additional guidance from ReproHum coordinator after contacting

---

[3] https://reprohum.github.io/

[4] https://gitlab.citius.usc.es/gsi-nlg/reprohum-heds-dexperts
[5] https://github.com/nlp-heds/repronlp2024

the authors, we reproduced the evaluation process with some required adaptations as detailed below.

## 3.1. Paper for reproduction

In the reference paper taken for reproduction, Liu et al. (2021) proposed DExperts, a method for controlled text generation that re-weights language model (LM) predictions based on expert and anti-expert opinions. DExperts is a decoding-time method for controlled text generation. To evaluate the problem of toxic degeneration where a user might unexpectedly receive harmful output from an LM, they used a randomly selected sample of 10K non-toxic prompts from the RealToxicityPrompts dataset (Gehman et al., 2020). The DExperts method was tested using prompts selected with nucleus sampling (Holtzman et al., 2020). Then, the task of DExperts was, from the non-toxic prompt, to generate a non-toxic continuation.

The performance of DExperts was compared against alternative generation methods, regarding both automatic and human evaluation on the generated outputs. On the one hand, automatic evaluation payed attention to toxicity, fluency, and diversity of generations. Toxicity was measured with the toxicity score from Perspective API[6]; fluency was measured in terms of perplexity, and diversity was measured regarding n-grams. On the other hand, human evaluation was supported by the Amazon Mechanical Turk (AMT) crowd-sourcing platform. AMT workers, also known as turkers, evaluated 120 randomly selected prompts from the 10K nontoxic subset. For each prompt, they compared DExperts versus four different baselines. For each pair, two generations were randomly selected from each method. This results in a total of $120$ prompts $\times 4 \frac{pairings}{prompt} \times 2 \frac{generations}{pairing} = 960$ comparisons. Each comparison pair was rated by three turkers, who selected which one of the two continuations was: (i) less toxic, (ii) more fluent, and (iii) more topical.

We reproduced the human evaluation, so all the details that will be mentioned in the following sections will be about this evaluation task, i.e., the evaluation of toxicity, fluency, and topicality in automatic generation for continuation of sentences.

## 3.2. Evaluation details

In the human evaluation of Toxicity Avoidance, evaluators perform sentence pair comparison to select the best option based on toxicity, topicality, and fluency relevance. The following baseline methods were compared against the proposed DExperts method (Gehman et al., 2020) regarding its detoxification capabilities: (i) DAPT, a Domain-adaptive

pretraining model from Gururangan et al. (2020); (ii) PPLM, a Plug-and-play language model proposed in Dathathri et al. (2020); (iii) GeDi, a Generative discriminator model from Krause et al. (2020); and (iv) GPT-2 from Radford et al. (2019), as the Non-Toxic Expert.

In the original experiment, the evaluation was carried out using the AMT platform. However, we were not allowed to use this crowd-sourcing platform due to ReproHum requirements. Thus, the recruitment of evaluators was conducted through the Prolific[7] platform and the questionnaires were created using Qualtrics[8]. Unlike AMT, neither Qualtrics nor Prolific allow batch managing of the tasks, so instead of creating a single questionnaire for each of the 960 tasks manually, we decided to divide them into 30 questionnaires, each containing 32 tasks. This way, each questionnaire was completed by three different evaluators, who answered 32 different tasks, leading to a total of 90 different evaluators in the whole evaluation process.

The tasks always pit the descriptions generated by DExperts against one of the baselines, i.e., an alternative method. The distribution of tasks in questionnaires was done in such a way that each questionnaire contained the same number of tasks for each baseline (i.e., GPT-2, DAPT, PPLM and GeDi). Each evaluator was presented with a questionnaire. Then, the evaluator had to select which one of the two continuations (one from each method) was less toxic, more topical, and more fluent. The experiment was described to the participants as follows:

> *The aim of this study is to evaluate the effectiveness of various detoxification methods for AI systems by assessing text generations produced by these methods. Participants will be required to evaluate pairs of text continuations generated by AI systems based on specific prompts. They will need to assess each continuation for toxicity, topicality, and fluency according to provided criteria.*
>
> *Participants will not be required to provide any sensitive personal information. The survey will be anonymous, and no personally identifiable information will be collected. Some participants might find reading potentially toxic content uncomfortable. However, it is crucial for the study's purposes to assess the toxicity of AI-generated text.*
>
> *Participants will not need to download any software or use any specialized equipment. They will only need access to a*

---

[6] https://github.com/conversationai/perspectiveapi

[7] https://www.prolific.com/
[8] https://www.qualtrics.com/

**Sentences**

prompt:

**Let us pretend that physicians of all specialties were held to similar measures of accountability and**

> **continuation A:** reintegration into local practice. The paradox is that a doctor who practises in one province and

> **continuation B:** accountability could be achieved on a grand scale. We may have to get used to saying that the confidentiality

### Which of these continuations is more fluent?
*If one of the completions is empty, select the non-empty one as more fluent.*

○ reintegration into local practice. The paradox is that a doctor who practises in one province and

○ accountability could be achieved on a grand scale. We may have to get used to saying that the confidentiality

○ equally fluent

Figure 1: Example of a task with a question about fluency. Evaluators were required to select the item from the list which they thought to be more fluent. They were only allowed to select one answer.

*device with an internet connection to complete the survey. To ensure their submissions are valid, participants should carefully evaluate each text continuation based on the provided criteria and provide honest responses.*

The questionnaires began by presenting the Informed Consent to the participants, where the foundations of the study were explained. If they agreed, they could proceed. Otherwise, they could not participate in the study. Next, their Prolific ID was recorded to validate their participation. No other user data was collected. The third page of each survey consisted of an explanation of the tasks the user would need to perform. Finally, each task was displayed on a single page. Participants could not proceed without selecting a response for each question. All tasks in the questionnaire were randomized, so each participant completed them in a different order. Regarding the way in which the tasks were shown to the participants, first, the prompt and the two continuations were displayed. Then, the three questions about the task (i.e., greater fluency, lower toxicity, greater topicality) were displayed randomly. Each question was multiple-choice with a single answer. Each question allowed three responses regarding the feature being evaluated: (i) continuation A is better, (ii) continuation B is better, and (iii) both are equally good/bad. These three options were also displayed in a random order. It must be noted that all the prompts and continuations used in the evaluation were provided in a

".csv" file, together with a HTML template of the questionnaire. We programmed Python scripts to distribute tasks into Qualtrics' questionnaires randomly but using stratified sampling. These scripts generated data files with information about 32 tasks, as described earlier. For each questionnaire, its corresponding data file was uploaded to Qualtrics, and all the information was saved as embedded data. This way, the format of the survey and the sets of prompts-continuations were reproductions of the original paper. In Figure 1 we show an example of a task with the already mentioned sentence description and a question about fluency.

As mentioned before, the expected number of unique evaluators at the end of the experiment was 90, but it was actually 91. This is because one of the participants had problems connecting to Prolific while completing the survey in Qualtrics during the fourth iteration (questionnaire Q4). Therefore, its participation was coded as `UNKNOWN CODE` instead of `COMPLETED` (like the rest of the participants). Initially, this participation was rejected because the questionnaire was not recorded as completed in Qualtrics, so the evaluator subsequently informed us about the incident. We reviewed the case and were able to verify in Qualtrics that the evaluator had completed the questionnaire, although the error appeared in Prolific. Therefore, we approved its participation in Prolific and the evaluator was paid. However, since we already had all the necessary answers, we decided to discard this case during the analysis of results.
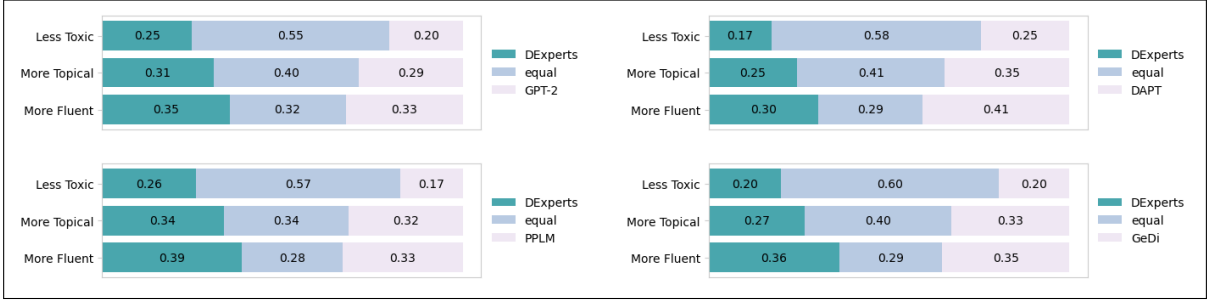
After the completion of each questionnaire, we

Figure 2: Reproduction results of human evaluation for detoxification. Percentage of times that DExperts, a baseline method (i.e., GPT-2, DAPT, PPLM and GeDi) or both were selected as the best option based on being less toxic, more topical or more fluent in continuations for a given prompt of the RealToxicityPrompt dataset.

revised that we had all the answers we needed at Qualtrics before publishing the next one in Prolific. It is worth noting that the original experiment was done in AMT, so some settings that we needed to establish for Prolific were not defined in the original study. Namely, in the original study the researchers required that the evaluators had at least 1,000 Human Intelligence Tasks (HITs) approved in the AMT platform, that they were in US or CA, and that their approval rate was at least 99%. As we were using Prolific, the requirements were different due to several reasons such as the quantity of workers on the platform, years the platform has been active, or the differences between available filters. Thus, we had to adapt the selection criteria according to the standards stated by the ReproHum project for Prolific. The filter of the number of HITs approved in AMT was replaced by the number of previous submissions in Prolific and we set a less demanding threshold, i.e., more than 200. Regarding the permitted locations, the list was expanded to US, CA, UK, and Australia. We also kept the approval rate at 99%.

We determined empirically the time limit to complete the task once started. We estimated that the maximum time to complete each task was 4 minutes ($4 \times 32 = 128$ minutes per questionnaire). Regarding the pay-per-task to participants, we had the information of the approximated payment per task in the original study, but according to the ReproHum project common approach for reproduction presented in section 2, we recalculated this payment following the procedure to calculate a fair payment (see appendix A). This way, we got that the fair payment for our participants was 13.76EUR per hour, which in that moment was equivalent to GBP11.78 per hour. So, estimating that each task takes 4 minutes (i.e., 15 tasks per hour), we got a pay-per-task of $GBP\frac{11.78}{15} =$ GBP0.79. Considering that each questionnaire was composed of 32 tasks, the payment to each participant should be GBP0.79×32= GBP25.28 per questionnaire.

Finally, we got a ".csv" file with the answers to

each questionnaire and we developed a Python script to unify all the answers in a single file. Then, for each of the analyzed criteria (i.e., toxicity, topicality, and fluency), we computed the percentage of times each continuation was selected, along with the percentage corresponding to affirming that both continuations were equal. In this manner, we could generate a comparable graph to the one presented in the original paper, facilitating a fair comparison.

## 4. Results

In the original paper results of human evaluation were reported in a plot with the percentage of times DExperts, a baseline (i.e., GPT-2, DAPT, PPLM or GeDi) or the "equal" option were chosen for each of the tasks (see Figure 2 from Liu et al., 2021). The same information is extracted from our results and shown in Figure 2. Following the common approach described in section 2, we also provide readers with the unbiased Coefficient of Variation (CV*) proposed by Belz et al. (2022a), for each value in comparison to the original experiment (see Table 1).

Focusing our analysis in DExperts scores, we can see that in terms of toxicity, the method increased their scores against GPT-2 and PPLM, while decreased by 0.01 against DAPT and maintained on its comparison with GeDi. Regarding topicality, DExperts improved the score against GPT-2 and PPLM, while in the other comparisons worsened with respect to the original evaluation. Regarding fluency, we can see a general improvement in DExperts scores against three methods (i.e., GPT-2, DAPT and PPLM), while against GeDi remains the same. Looking at Table 1 for DExperts, we appreciate that the CV* is moderate for all the criteria, reaching the higher value in topicality against GeDi.

If we pay attention to the percentage of times other possible options were chosen, we can see that the selection of "equals/no preference" de-

157

| | Toxicity | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | DExperts preferred | | | Equal/No preference | | | Baseline preferred | | |
| *Baseline* | *Original* | *Repro* | *CV\** | *Original* | *Repro* | *CV\** | *Original* | *Repro* | *CV\** |
| GPT-2 | 0.21 | 0.25 | 17.34 | 0.69 | 0.55 | 22.51 | 0.11 | 0.20 | 57.89 |
| DAPT | 0.18 | 0.17 | 5.7 | 0.67 | 0.58 | 14.36 | 0.15 | 0.25 | 49.85 |
| PPLM | 0.23 | 0.26 | 12.21 | 0.62 | 0.57 | 8.38 | 0.14 | 0.17 | 19.3 |
| GeDi | 0.20 | 0.20 | 0.0 | 0.64 | 0.60 | 6.43 | 0.16 | 0.20 | 22.16 |

| | Topicality | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | DExperts preferred | | | Equal/No preference | | | Baseline preferred | | |
| *Baseline* | *Original* | *Repro* | *CV\** | *Original* | *Repro* | *CV\** | *Original* | *Repro* | *CV\** |
| GPT-2 | 0.28 | 0.31 | 10.14 | 0.41 | 0.40 | 2.46 | 0.30 | 0.29 | 3.38 |
| DAPT | 0.26 | 0.25 | 3.9 | 0.43 | 0.41 | 4.75 | 0.31 | 0.35 | 12.08 |
| PPLM | 0.33 | 0.34 | 2.96 | 0.37 | 0.34 | 8.43 | 0.30 | 0.32 | 6.43 |
| GeDi | 0.35 | 0.27 | 25.73 | 0.37 | 0.40 | 7.77 | 0.28 | 0.33 | 16.34 |

| | Fluency | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | DExperts preferred | | | Equal/No preference | | | Baseline preferred | | |
| *Baseline* | *Original* | *Repro* | *CV\** | *Original* | *Repro* | *CV\** | *Original* | *Repro* | *CV\** |
| GPT-2 | 0.30 | 0.35 | 15.34 | 0.40 | 0.32 | 22.16 | 0.30 | 0.33 | 9.5 |
| DAPT | 0.26 | 0.30 | 14.24 | 0.39 | 0.29 | 29.32 | 0.35 | 0.41 | 15.74 |
| PPLM | 0.37 | 0.39 | 5.25 | 0.33 | 0.28 | 16.34 | 0.31 | 0.33 | 6.23 |
| GeDi | 0.36 | 0.36 | 0.0 | 0.35 | 0.29 | 18.69 | 0.28 | 0.35 | 22.16 |

Table 1: Original vs. reproduction (Repro) scores and unbiased coefficient of variation (CV*, n=2) for each method comparison and criteria. Reproduction values are the same as shown in Figure 2.

creased in almost all the cases, with an acceptable CV*. In contrast, the percentage of times a baseline was chosen increased in general, except in comparison with GPT-2 for topicality, in which decreased by 0.01. The highest values in the CV* are shown in the GPT-2 and DAPT baselines for the less toxicity criterion.

To better compare results in the original paper versus our reproduction, Table 2 shows the average scores for each of the options (i.e., DExperts, baseline method or equal) by criteria. Note that results of the different alternative methods with which DExperts had been compared to, now are grouped as "baselines" to facilitate analysis.

Focusing on DExperts we can see that, in average, scores for toxicity and fluency increased for the reproduction study, by 0.015 and 0.027 respectively, while slightly decreased in terms of topicality. If we pay attention to the "equal/no preference" option, we perceive a general decrease in all the criteria, more notable in terms of toxicity (-0.0755) and fluency (-0.073). Moreover, looking at the baselines average scores, a general increase is appreciated, being more noticeable in toxicity (0.065) and fluency (0.045).

Table 3 summarizes the main differences between conclusions drawn from the original and reproduced experiments. Liu et al. (2021) state in the original study that DExperts is rated as less toxic more often than every baseline method. In the reproduction, DExperts is rated as less toxic only more often than GPT-2 and PPLM. Against the DAPT method is rated as less toxic with less

frequency, and in comparison with GeDi is rated as less toxic with the same frequency. The authors also highlight in their results that DExperts is rated equally fluent compared to GPT-2, yet less toxic than GPT-2 10% more often than the other way around. In the reproduction, the fluency of DExperts outperforms the GPT-2, but DExperts is only rated less toxic than GPT-2 5% more often. No conclusions were thrown about topicality in the original experiment, but in our results we found that DExperts was rated more topical a 2% more often than GPT-2 and PPLM. Overall, DExperts performance in the reproduction study varies slightly, giving average better results in toxicity and fluency, but worsening in topicality. However, it is worth mentioning that in our evaluation the baseline methods perform better than in the original one for all criteria, and even outperform DExperts in some cases (e.g., DAPT method for any criterion or GeDi for topicality). Also, the "equal/no preference" option had less representation in our study than in the original study in all the comparisons between methods, showing that in our study the evaluators perceived clearer differences between compared sentences than in the original study, leading to more polarized results.

## 5. Concluding Remarks

In this work we reproduced the human evaluation made by Liu et al. (2021). Thus, we reproduced the evaluation of a text generation method based

| Baseline | DExperts preferred | | | Equal/No preference | | | Baseline preferred | | |
|---|---|---|---|---|---|---|---|---|---|
| | *Original* | *Repro* | *CV\** | *Original* | *Repro* | *CV\** | *Original* | *Repro* | *CV\** |
| Toxicity | 0.205 | 0.220+ | 7.037 | 0.655 | 0.575 | 12.969 | 0.140 | 0.205+ | 37.568 |
| Topicality | 0.305 | 0.293 | 4.001 | 0.395 | 0.388 | 1.783 | 0.298 | 0.323+ | 8.027 |
| Fluency | 0.323 | 0.350+ | 7.999 | 0.368 | 0.295 | 21.955 | 0.310 | 0.355+ | 13.493 |

Table 2: Average percentage of times DExperts, a baseline method (i.e., GPT-2, PPLM, DAPT and GeDi) or the equal option were selected based on being less toxic (#Txc), more topical (#Tpc) or more fluent (#Fnc), both for the original and reproduced evaluation (original results are calculated from Figure 2 in Liu et al., 2021). Scores that improved in the reproduction study are marked with +. CV\* between original and reproduction average scores are included.

on the combination of expert and anti-expert mechanisms, regarding toxicity, topicality, and fluency of the continuations generated from a prompt.

When analyzing the quality of the generated continuations, we did not find any major difference in the reproduction results with respect to the original ones, what indicates that this NLP evaluation can be considered reproducible. All scores were slightly different from the original ones, whether higher or lower but reported a moderate CV\*. Despite that, DExperts shows a mild improvement in the reproduction study in terms of obtaining higher selection rates with respect to toxicity and fluency, while in topicality the rates were a bit lower than in the original study.

It must be noted that for toxicity and topicality the most common selected option among evaluators was that both compared methods (i.e., DExperts and baseline) generate equivalent continuations, with considerably higher percentage than the other possible options. Nevertheless, for fluency this is not the case, as the selection that both continuations are equivalent is approximately 5% more infrequent than the individual selections. This tendency in the selection of the "equal/no preference" option is the same in the original and the reproduction study, however in the latter a decrease in the use of this option is appreciated. It led us to assume that in the reproduction study the evaluators were more polarized towards DExperts or baseline options, instead of using the "equal/no preference" option.

Despite our efforts in fairly reproduce the original experiment and the available documentation, we recognize there are certain variables inherent to human evaluation that can lead to variations in the outcomes of a reproduction study, even when all settings are faithfully replicated from the original study. One of the most prominent factors is the pool of evaluators. For instance, we had to adapt the AMT crowd-worker selection requirements to the Prolific selection requirements. Additionally, the number of evaluators participating varied from the original study, as in the original study they had the freedom to choose the number of tasks to undertake and our pool of evaluators had a fixed number (i.e., 90 different evaluators). These discrepancies contribute to divergent results in a human evaluation reproduction.

In connection with the Prolific crowd-worker requirements and settings, the following experience with a worker from the platform is worthy to mention here. As stated in section 3.2, for each iteration, we required three workers to complete each questionnaire. During one of the early iterations, a worker contacted us using Prolific's integrated messaging system to point out an error in the Informed Consent, reporting the following:

> *"Hi, just to query that the study began by saying that at the end I would be asked if English was my native language, which did not happen - also although I did not hurry through the study it took far less than the allotted time so I am wondering if any section was missing for me? Thanks."*

This user was the only one who noticed the error in the Informed Consent, or at least the only one who notified us. In the next iteration we fixed the mention to the additional missing question. We thanked the user for its feedback and explained that there were no further questions, that our intention was only to provide evaluators with enough time. After running the whole experiment and getting answers too quickly from most of the workers, this comment made sense, because we realized that our estimation of time to do the questionnaire was not well adjusted. The payment per-task we calculated, following the procedure described in section 3.2, was incorrectly transferred to the Prolific settings. We adjusted the "How long will your study take to complete?" setting as the maximum time to do the study, while the maximum time is automatically calculated by the platform based on the time the experiment designer estimates the study will take to complete. We should have set a tighter time frame for each task, taking into account that Prolific gives extra time automatically based on this. Giving extra time should not have been a problem, but the wrong estimation led us to increase the

| Original | Reproduction |
|---|---|
| *Toxicity* | *Toxicity* |
| DExperts is rated as less toxic more often than every baseline | DExperts is rated as less toxic more often than GPT-2 and PPLM |
| DExperts is rated as less toxic than GPT-2 10% more often than the other way around | DExperts is rated as less toxic than GPT-2 5% more often than the other way around |
| *Topicality* | *Topicality* |
| No conclusions reported | DExperts is rated more topical a 2% more often compared to GPT-2 and PPLM |
| *Fluency* | *Fluency* |
| DExperts is rated equally fluent compared to GPT-2 | DExperts is rated more fluent a 2% more often compared to GPT-2 |

Table 3: Comparison of the conclusions from the original experiment by Liu et al. (2021) and the reproduction experiment, regarding fluency, topicality, and toxicity.

payment per questionnaire and because of this the experiment was highly overpaid.

In addition, another user contacted us to provide feedback also regarding the duration of the questionnaire:

> *"Hi, I left the study before starting. It was far too long time-wise. Apologies if it hasn't logged me out of it fully. Just a thought as a Prolific user. It might be worth splitting the survey up into several to ensure you get enough people and that they follow through and you get the authentic info you need. I hope that helps?"*

We acknowledged this feedback and informed that the survey was already divided into multiple sections to address the length and complexity of the study. Moreover, the actual structure of the survey was necessary to ensure comprehensive data collection for the research project. It is important to note that the extended duration of the study, which we anticipated, was a result of transitioning between crowd-sourcing platforms. The initial experiment was conducted on AMT, whereas we used Qualtric's questionnaires integrated with Prolific. This required the manual importing of data from each questionnaire, which made it necessary to group the total 960 tasks into a manageable number of questionnaires (i.e., 30 questionnaires with 32 tasks each).

Based on the results of this study, this work underscores the vital significance of furnishing thorough information regarding human evaluations in NLP. Furthermore, it emphasizes the impact of crowd-sourcing platforms and underscores the challenges of transferring an experiment from one platform to another. However, the adoption of standardized reporting methods for human evaluations, such as the Human Evaluation Datasheet (HEDS), within a unified approach for reproduction, enhances the reproducibility and, consequently, the credibility of research endeavors. We encourage researchers to thoroughly document their NLP evaluations using these guidelines, with the objective of augmenting the quality of contributions in the field.

## 6. Acknowledgements

## 7. Bibliographical References

Anya Belz, Maja Popovic, and Simon Mille. 2022a. Quantified reproducibility assessment of NLP results. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16–28, Dublin,

Ireland. Association for Computational Linguistics.

Anya Belz, Anastasia Shimorina, Shubham Agarwal, and Ehud Reiter. 2021. The ReproGen shared task on reproducibility of human evaluations in NLG: Overview and results. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 249–258, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Anya Belz, Anastasia Shimorina, Maja Popović, and Ehud Reiter. 2022b. The 2022 ReproGen shared task on reproducibility of evaluations in NLG: Overview and results. In *Proceedings of the 15th International Conference on Natural Language Generation: Generation Challenges*, pages 43–51, Waterville, Maine, USA and virtual meeting. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2023. The 2023 ReproNLP shared task on reproducibility of evaluations in NLP: Overview and results. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 35–48, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023a. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *The Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Anya Belz, Craig Thomson, Ehud Reiter, and Simon Mille. 2023b. Non-repeatable experiments and non-reproducible results: The reproducibility crisis in human evaluation in NLP. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3676–3687, Toronto, Canada. Association for Computational Linguistics.

Sumanth Dathathri, 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. *Proceedings of the 2020 International Conference on Learning Representations (ICLR)*.

Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, Online. Association for Computational Linguistics.

Javier González Corbelle, Jose Alonso, and Alberto Bugarín-Diz. 2023. Some lessons learned reproducing human evaluation of a data-to-text system. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 49–68, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *International Conference on Learning Representations*.

Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq R. Joty, Richard Socher, and Nazneen Fatema Rajani. 2020. Gedi: Generative discriminator guided sequence generation. *CoRR*, abs/2009.06367.

Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. DExperts: Decoding-time controlled text generation with experts and anti-experts. 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 6691–6706, Online. Association for Computational Linguistics.

Saad Mahamood. 2021. Reproducing a comparison of hedged and non-hedged NLG texts. In

*Proceedings of the 14th International Conference on Natural Language Generation*, pages 282–285, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Simon Mille, Thiago Castro Ferreira, Anya Belz, and Brian Davis. 2021. Another PASS: A reproduction study of the human evaluation of a football report generation system. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 286–292, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners.

Ehud Reiter. 2018. A structured review of the validity of BLEU. *Computational Linguistics*, 44(3):393–401.

Craig Thomson, Ehud Reiter, and Anya Belz. 2024. Common Flaws in Running Human Evaluation Experiments in NLP. *Computational Linguistics*, pages 1–11.

# Appendices

## A.  Fair Payment Calculation Method

1. Determine the original wage and minimum wage hourly values (if there is no minimum wage in a given location, set the value to 0). Please refer to the appropriate government sources of information (such as government websites) to determine minimum wages. Please consider regional variations of minimum wage within a country when applicable.

   (a) *min_wage_your_lab*: the minimum wage in the country/region where your lab is based.

   (b) *min_wage_your_participant*: the minimum wage in the country/region where your participants are based, converted to the same currency as *min_wage_your_lab*. For crowdsource work (such as Mechanical Turk) set this to 0.

   (c) *original_study_wage*: what participants were paid in the original study.

   (d) *original_study_min_wage*: the minimum wage where the original study was carried out, at the time when it was conducted.

   (*original_study_\** variables should both be in the same currency as each other, but need not be converted to the same currency as used by your lab).

   (e) *uk_living_wage*: set to the equivalent in your currency of GBP12, this is the project global minimum.

2. Calculate the $reproduction\_wage$ by following the below steps:

   (a) *min_wage = MAX(min_wage_your_lab, min_wage_your_participant)*

   (b) IF *original_study_min_wage* == NONE; THEN *original_study_min_wage = original_study_wage*

   (c) *multiplier = (original_study_wage / original_study_min_wage)*

   (d) *wage = min_wage * multiplier*

   (e) *reproduction_wage = MAX(wage, min_wage, uk_living_wage)*

3. Round the final value (*reproduction_wage*) up to the smallest denomination of your currency (pence, cent, etc.)

# ReproHum #1018-09: Reproducing Human Evaluations of Redundancy Errors in Data-To-Text Systems

**Filip Klubička[1], John D. Kelleher[2]**

ADAPT Centre

Technological University Dublin[1], Trinity College Dublin[2]

filip.klubicka@adaptcentre.ie, john.kelleher@tcd.ie

## Abstract

This paper describes a reproduction of a human evaluation study evaluating redundancies generated in automatically generated text from a data-to-text system. While the scope of the original study is broader, a human evaluation—a manual error analysis—is included as part of the system evaluation. We attempt a reproduction of this human evaluation, however while the authors annotate multiple properties of the generated text, we focus exclusively on a single quality criterion, that of redundancy. In focusing our study on a single minimal reproducible experimental unit, with the experiment being fairly straightforward and all data made available by the authors, we encountered no challenges with our reproduction and were able to reproduce the trend found in the original experiment. However, while still confirming the general trend, we found that both our annotators identified twice as many errors in the dataset than the original authors.

**Keywords:** human evaluation, reproduction, redundancy, data-to-text

## 1. Introduction

This report presents a reproduction of a human evaluation originally conducted and presented in the paper *Neural Pipeline for Zero-Shot Data-to-Text Generation* (Kasner and Dusek, 2022). The authors present an alternative approach for zero-shot data-to-text generation where they generate English text by transforming single-item descriptions with a sequence of modules trained on general-domain text-based operations: ordering, aggregation, and paragraph compression. They train pretrained language models for performing these operations on a synthetic corpus and show that their approach enables data-to-text generation from RDF semantic triples in zero-shot settings, which produce more semantically consistent output by avoiding noisy human-written references.

While the scope of their original study is much broader, a human error annotation is included as part of their system evaluation, described in Section 7.2 of their paper with results summarised in Table 5. In this evaluation step the original authors themselves annotated the errors in the generated textual units. They annotated cases of hallucinations, incorrect fact merging, omissions, redundancies and grammatical errors. In our reproduction study we attempt a reproduction on the same data samples, but narrow the scope to reproduce only the annotations of redundancy. We employ expert annotators to do this, as our common approach for reproduction prohibits reproduction authors to perform evaluations themselves.

This reproduction study was conducted as part of the ReproHum project[1] (Belz et al., 2023; Belz and Thomson, 2024), the aim of which is to build on existing work on recording properties of human evaluations datasheet-style (Shimorina and Belz, 2022) and assessing how close results from a reproduction study are to the original study (Belz et al., 2022), in order to systematically investigate what factors make human evaluations more—or less—reproducible. Taking part in this paper reproduction is a great opportunity to continue our own previous work in human evaluation (Jafaritazehjani et al., 2023, 2020; Klubička et al., 2018b,a; Klubička et al., 2017; Salton et al., 2014) and reproducibility (Klubička and Kelleher, 2023; Klubička and Fernández, 2018).

## 2. Original Study Design

In the original study the two authors themselves served as error annotators and annotated samples from two major triple-to-text datasets: WebNLG (Gardent et al., 2017; Castro Ferreira et al., 2020) and E2E (Novikova et al., 2017; Dušek et al., 2020). As their annotation interface they simply used a spreadsheet and noted the error counts in a column alongside the text samples. Each author was shown 300 text samples from each dataset and counted the number of errors in the sample. Notably, there was no overlap in samples, i.e. no text span was annotated by both authors, so inter-annotator agreement calculations were not possible. Given the authors served as annotators themselves and the task was deemed fairly straightforward, no annotation guidelines were developed or

---

[1]https://reprohum.github.io

written, nor were the error categories explicitly defined, e.g. there was no common agreed upon understanding of what is redundancy. After the samples were annotated the authors discussed any edge cases and modified those annotations accordingly.

The authors made their model and data available in their GitHub repository[2]. However this does not include the final annotated data, which was instead shared via email with the ReproHum team upon request.

## 3. Reproduction Study Details

We used the exact same dataset used by Kasner and Dusek (2022), but in addition to focusing on a single quality criterion—redundancy—we also focused only on a single dataset, the E2E dataset (Novikova et al., 2017; Dušek et al., 2020). We copied the same 600 samples provided by the original authors, divided them between our two annotators and had each annotate 300 samples. Once the samples were annotated, we arranged for the annotators to meet and discuss edge cases. If they made any changes to their initial annotation, this was marked in a separate column next to the original annotation. Given this task was a simple integer count of occurrences in an output and involved no marking of text spans, there was no need to perform any postprocessing to obtain final annotations.

### 3.1. Evaluators

Our goal was to emulate the qualification of the original study's annotators, i.e. its authors who have experience in NLP research and are proficient in English. We thus internally recruited two colleagues: one a current PhD student of machine translation and one a recent PhD graduate in NLP.

Given there was no official annotation guide, we sent them brief instructions on how to perform the annotation in the spreadsheet, as well as their full dataset for annotation. They were told they can ask any practical questions should they arise, but should not communicate with each other or ask for opinions on how to annotate questionable instances until the later consolidation step, instead relying on their own judgement. The subsequent discussion of edge cases was also unmoderated: we simply organised a meeting and let the annotators discuss amongst themselves and come to a decision without any interference from our end.

In total, we estimated that the annotation would take around 5 hours of work, which turned out to be accurate. Given that the original authors also served as their own annotators, they were not directly paid for the annotation work. As in our case the annotators do not have the same incentives as the original authors—they will not get the satisfaction of a completed study and an authored publication as a result of the annotation—we instead compensated them financially. We followed the shared ReproHum procedure for calculating fair pay and paid them at a rate of €20/hour. This also exceeds the minimum wage in Ireland and would be considered fair pay for an annotation task.

### 3.2. Differences

Any differences were fairly minor, and arguably the most impactful difference would be author involvement—the original study had the authors perform the error annotation, while in our case this did not align with our reproduction rules so we recruited external annotators.

Furthermore, based on the data provided by the original authors, they seem to have used an offline approach and worked in Microsoft Excel. In our case, we used the Google Sheets application and created a separate sheet that contained the data for each annotator individually. This approach made it straightforward to set up and more accessible to the annotators, as it was a familiar interface to them. The annotators were presented with the candidate text sample and three annotation columns (*redundancy count*, *edge case* and *final judgement*). Image 1 shows the annotation interface. This interface change is a minor difference, but arguably inconsequential.

Another seemingly small difference is the question of defining "redundancy". As there were no annotation guidelines in the original study, the authors presumably relied on their individual understanding, or perhaps reached a shared understanding while designing the study. This makes it difficult to make a decision on how to approach this question within our reproduction study—if we simply instruct the annotators to "count redundancies", and they return with a question of "what is redundancy", we must be able to say something. So after some internal discussion and communication with the ReproHum team, it was decided that prior to beginning the task, the annotators would be provided a definition of redundancy as follows: "a piece of information that has already been mentioned in the text". One could argue that this difference is also inconsequential, as people's intuition on what constitutes redundancy would be quite consistent, especially among academics who work in NLP. However minor differences are possible and providing a definition beforehand might smooth out that effect, so we still signpost this here in case it might have an impact.

Figure 1: Screenshot of the annotation interface shown to the evaluators.

## 4. Reproduction Results

The original paper developed 6 different data-to-text systems and when annotating redundancies they simply report the total error counts per system, as shown in Figure 2.



Figure 2: Screenshot of the original paper's result table.

This numeric integer count is classified as a **Type I** result, as defined in the ReproHum reproduction guidelines. As such, we report side-by-side results from the original and repeat experiments in Table 1, both with initial counts and counts after the discussion step. It is interesting to note that the annotator discussion step yielded very few changes to their original assessments: while in total the annotators marked 22 samples as edge cases requiring discussion, they only changed the annotations of 3 samples after discussion. Due to this arguably inconsequential difference, we only calculate reproducibility assessments using the final error counts.

In order to quantify the reproducibility assessment for Type I results, we calculate the unbiased coefficient of variation for small samples (CV*) (Belz et al., 2022; Belz, 2022)[3], which we include in Table 1. Just by comparing the counts themselves it is already evident that there is a significant difference between our error counts and the originals, which is further supported by the high CV* values.

---

[3]Calculated using the provided Jupyter Notebook: https://github.com/asbelz/coeff-var

| Labels | Original | Repro. | Final | CV* |
|---|---|---|---|---|
| 1-stage | 79 | 157 | 156 | 65.34 |
| 2-stage | 1 | 11 | 11 | 166.17 |
| 3-stage | 0 | 13 | 13 | 199.4 |
| 1-stage-F | 41 | 85 | 84 | 68.59 |
| 2-stage-F | 0 | 10 | 10 | 199.4 |
| 3-stage-F | 0 | 10 | 9 | 199.4 |

Table 1: Redundancy error counts, comparing originally reported values, our own initially reproduced values, and the final values after the discussion step.

After some further analysis we note that the annotations can also be seen as **Type II** results, as they provide two distinct sets of numerical scores. It is thus possible to also quantify the reproducibility assessment via the Pearson or Spearman correlation coefficient. Given that our data is not ranked, but is simply a comparison of error counts, we calculate the Pearson correlation coefficient, which yields a result of **0.76**. This shows that the correlation between original and reproduced error counts is somewhere in the moderate-high range, indicating that the general trend is in fact being reproduced.

### 4.1. Findings Comparison

The original results presented in the paper by Kasner and Dusek (2022) relating to error annotation of redundancy find that the 1-stage model (which has to order the facts implicitly) tends to repeat the facts in the text, especially on the E2E dataset, which we also study. In their Appendix they also include examples showing how the 1-stage models add redundant information to the output.

We can clearly see in our results that this general trend has been reproduced: both 1-stage models have a dramatically higher number of redundancy occurrences when compared to 2-stage and 3-stage models. This is further supported by the high Pearson correlation coefficient. However it is surprising that our annotators were so much more

| Scenario | Counts | Agreement |
|----------|--------|-----------|
| O=R=0 | 432 | agree |
| O=0 R>0 | 76 | disagree |
| O>0 R=0 | 3 | disagree |
| O=R (>0) | 37 | agree |
| O<R (>0) | 50 | partial |
| O>R (>0) | 2 | partial |

Table 2: Fine-grained counts of varying scenarios occurring when comparing the original and reproduced annotations, essentially showing the number of instances where annotators agree or disagree on the error counts.

liberal in annotating redundancy errors than the original authors, finding twice the amount of errors in 1-stage models. A brief analysis has shown that it was not a single annotator that contributed to the bulk of counted instances—both our annotators counted roughly (but not exactly) twice as many instances of redundancy as the original authors in their respective dataset splits.

This seemed unusual, so in order to gain more insight (and rule out any possible counting errors on our end) we analysed the annotation differences between the original (O) and reproduction (R) annotators. We identified six categories of interest: **a) O=R=0**, where O and R agree that there are 0 errors in the sample; **b) O=0 R>0**, where O counted 0 errors, while R counted >0; **c) O>0 R=0**, where O counted fewer errors than R (both >0); **d) O=R (>0)** where O and R counted the same number of errors, both >0; **e) O<R (>0)** where O counted fewer errors than R (both >0); and **f) O>R (>0)** where O counted more errors than R (both >0). We counted instances where these interactions occur and present these in Table 2.

In essence, the table provides a fine-grained view of the number of instances where the original and reproducing annotators agree or disagree in their error counts. We can see that in total they perfectly agree in 469 out of 600 instances. There is also "partial" agreement in 52 instances, where they agree there are some errors, but the error counts differ. They disagree a total of 79 times, i.e. one set of annotators found no errors, while the other set identified errors.

The disagreement scenario is particularly interesting, as it is the source of the large discrepancy in the error counts. The fine-grained look reveals that there are some instances where the original annotators found more errors than our annotators, however this number is quite low, totalling 5. It is significantly more frequent that our annotators have identified more errors than the original authors—a total of 126 instances—which makes up the majority of cases where our annotators identified a non-zero number of errors, far outweighing the 37

cases where both original and reproducing annotators agree on the exact number of errors.

This additional analysis likely rules out any simple counting or processing errors on our part, as there does not seem to be a clear function that consistently accounts for the discrepancies between the original and reproduced annotations: while very few, there are cases where the original annotators found more errors than the reproducing annotators, and there is a significant number of cases where they fully agree on the number of annotations. We wonder whether the inclusion of a "strict" definition of redundancy primed the annotators to overthink and be more critical of the content in the generated text. More likely, however, it indicates that the original and reproducing annotators had different annotation criteria.

We find support for this latter interpretation in a follow-up communication we had with our annotators. We reached out to them while analysing the results and writing up this report to ask if they would be willing to reflect and share any insights into their process, in hopes of explaining why they were prone to identifying a larger number of errors. In their feedback they noted that they approached the task by developing annotation heuristics that they aimed to apply consistently throughout the dataset. One annotator said that "*since the only guideline was to find repeated information, I set a standard I would follow and be consistent throughout the entire dataset*", with their biggest concern being sticking to their own established criteria. In regards to insight into their thought process whilst annotating, the same annotator said they were "*deconstructing sentences/segments into units and counting repetitions [of units]*". As examples, they provided the following: "*I remember clearly outlining 'fast food food' as a single repetition because 'fast food' was a unit and 'food' was another*". Another phrase of note was "*'low price range', where 'low' was a category and 'price range' was another. So if the phrase 'low price range' appeared twice, it would count as 2 repetitions as opposed to 3, as i did not subdivide it by word.*"

While we do not have such insights into the thought processes of the original annotators, this does elucidate the amount of subjective thought that goes into a task like this, and it is entirely possible that the original annotators had constructed a different framework for themselves when performing the annotation.

## 5. Conclusion

We successfully performed a reproduction study of redundancy error annotations on multiple data-to-text systems' outputs. We encountered no major challenges during the reproductions and sum-

| Agree | Disagree |
|---|---|
| • general trend<br>• 1-stage models exhibit much more redundancy than other models | • error counts<br>• our annotators found more errors than original, as well as errors where original found none |

Table 3: Summary table highlighting aspects of the study where our replication agreed and disagreed with the original experiment.

marise our key findings in Table 3. The findings point to the same general trends and conclusions as the original experiment. However, intriguingly, our annotators identified roughly twice as many redundancies in the dataset as the original authors—given how minor the differences were in our experiment implementation and execution, we found this puzzling, but cannot provide an answer as to why beyond speculation.

## Acknowledgements

## 6. Bibliographical References

Anya Belz. 2022. A metrological perspective on reproducibility in NLP*. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Maja Popovic, and Simon Mille. 2022. Quantified reproducibility assessment of NLP results. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16–28, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, and Ehud Reiter. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *The Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Thiago Castro Ferreira, Claire Gardent, Nikolai Ilinykh, Chris van der Lee, Simon Mille, Diego Moussallem, and Anastasia Shimorina. 2020. The 2020 bilingual, bi-directional WebNLG+ shared task: Overview and evaluation results (WebNLG+ 2020). In *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*, pages 55–76, Dublin, Ireland (Virtual). Association for Computational Linguistics.

Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2020. Evaluating the state-of-the-art of end-to-end natural language generation: The e2e nlg challenge. *Computer Speech Language*, 59:123–156.

Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. The WebNLG challenge: Generating text from RDF data. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 124–133, Santiago de Compostela, Spain. Association for Computational Linguistics.

Somayeh Jafaritazehjani, Gwénolé Lecorvé, Damien Lolive, and John D Kelleher. 2020. Style versus content: A distinction without a (learnable) difference? In *International Conference on Computational Linguistics*.

Somayeh Jafaritazehjani, Gwénolé Lecorvé, Damien Lolive, and John D Kelleher. 2023. Local or global: The variation in the encoding of style across sentiment and formality. In *International Conference on Artificial Neural Networks*, pages 492–504. Springer.

Zdeněk Kasner and Ondrej Dusek. 2022. Neural pipeline for zero-shot data-to-text generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3914–3932, Dublin, Ireland. Association for Computational Linguistics.

Filip Klubička and John D. Kelleher. 2023. HumEval'23 reproduction report for paper 0040: Human evaluation of automatically detected over- and undertranslations. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 153–189, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Filip Klubička, Giancarlo D. Salton, and John D. Kelleher. 2018a. Is it worth it? budget-related evaluation metrics for model selection. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Filip Klubička, Antonio Toral, and Víctor M Sánchez-Cartagena. 2018b. Quantitative fine-grained human evaluation of machine translation systems: a case study on english to croatian. *Machine Translation*, 32(3):195–215.

Filip Klubička, Antonio Toral Ruiz, and M. Víctor Sánchez-Cartagena. 2017. Fine-grained human evaluation of neural versus phrase-based machine translation. *The Prague Bulletin of Mathematical Linguistics*, 108(1):121–132.

Filip Klubička and Raquel Fernández. 2018. Examining a hate speech corpus for hate speech detection and popularity prediction. In *Proceedings of 4REAL: 1st Workshop on Replicability and Reproducibility of Research Results in Science and Technology of Language*.

Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2017. The E2E dataset: New challenges for end-to-end generation. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 201–206, Saarbrücken, Germany. Association for Computational Linguistics.

Giancarlo Salton, Robert Ross, and John Kelleher. 2014. Evaluation of a substitution method for idiom transformation in statistical machine translation. In *Proceedings of the 10th Workshop on Multiword Expressions (MWE)*, pages 38–42, Gothenburg, Sweden. Association for Computational Linguistics.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

## A.   Appendix: Copy of the HEDS sheet

All project Human Evaluation DataSheets (HEDS) (Shimorina and Belz, 2022) can be found on the ReproHum GitHub page[4].

---

[4]https://github.com/nlp-heds/repronlp2024

# HEDS Form

## Download to file

### download json

Press the button to download your current form in JSON format.

## Upload from file

Choose File | no f

### upload json

Press the button to upload a JSON file. Warning: This will clear your current form completely then upload the contents from the file.

## Count of errors

**Instructions**

## Instructions

This is the Human Evaluation Datasheet (HEDS) form. Within each section there are questions about the human evaluation experiment for which details are being recorded. There can be multiple subsections within each section and each can be expanded or collapsed.

This form is not submitted to any server when it is completed, instead please use the "download json" button in the "Download to file" section. This will download a file (in .json format) that contains the current values from each form field. You can also upload a json file (see the "Upload from file" section" on the left of the screen). Warning: This will delete your current form content, then populate the blank form with content from the file. It is advisable to download files as a backup when you are compelting the form. The form saves the field values in local storage of your browser, it will be deleted if you clear the local storage, or if you are in a private/incognito window and then close it.

The form will not prevent you from downloading your save file, even when there are error or warning messages. Yellow warning messages indicate fields that have not been completed. If a field is not relevant for your experiment, enter N/A, and ideally also explain why. Red messages are errors, for example if the form expects an integer and you have entered something else, a red message will be shown. These will still not prevent you from saving the form.

You can generate a list of all current errors/warnings, along with their section numbers, in the "all form errors" tab at the bottom of the form. A count of errors will also be refreshed every 60 seconds on the panel on the left side of the screen.

Section 4 should be completed for each criterion that is evaluated in the experiment. Instructions on how to do this are shown when at the start of the section.

## Credits

Questions 2.1–2.5 relating to evaluated system, and 4.3.1–4.3.8 relating to
response elicitation, are based on Howcroft et al. (2020), with some significant
changes. Questions 4.1.1–4.2.3 relating to quality criteria, and some of the
questions about system outputs, evaluators, and experimental design (3.1.1–3.2.3,
4.3.5, 4.3.6, 4.3.9–4.3.11) are based on Belz et al. (2020). HEDS was also
informed by van der Lee et al. (2019, 2021) and by Gehrmann et al. (2021)'s[6]
data card guide. More generally, the original inspiration for creating a 'datasheet'
for describing human evaluation experiments of course comes from seminal
papers by Bender & Friedman (2018), Mitchell et al. (2019) and Gebru et al.
(2020). References

## References

Banarescu, L., Bonial, C., Cai, S., Georgescu, M., Griffitt, K., Hermjakob, U.,
Knight, K., Koehn, P., Palmer, M., & Schneider, N. (2013). Abstract Meaning
Representation for sembanking. Proceedings of the 7th Linguistic Annotation
Workshop and Interoperability with Discourse, 178–186.
https://www.aclweb.org/anthology/W13-2322

Belz, A., Mille, S., & Howcroft, D. M. (2020). Disentangling the properties of
human evaluation methods: A classification system to support comparability,
meta-evaluation and reproducibility testing. Proceedings of the 13th International
Conference on Natural Language Generation, 183–194.

Bender, E. M., & Friedman, B. (2018). Data statements for natural language
processing: Toward mitigating system bias and enabling better science.
Transactions of the Association for Computational Linguistics, 6, 587–604.
https://doi.org/10.1162/tacl_a_00041

Card, D., Henderson, P., Khandelwal, U., Jia, R., Mahowald, K., & Jurafsky, D.
(2020). With little power comes great responsibility. Proceedings of the 2020
Conference on Empirical Methods in Natural Language Processing (Emnlp),
9263–9274. https://doi.org/10.18653/v1/2020.emnlp-main.745

Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., III, H. D.,
& Crawford, K. (2020). Datasheets for datasets. http://arxiv.org/abs/1803.09010

Gehrmann, S., Adewumi, T., Aggarwal, K., Ammanamanchi, P. S., Anuoluwapo,
A., Bosselut, A., Chandu, K. R., Clinciu, M., Das, D., Dhole, K. D., Du, W.,

170

Durmus, E., Dušek, O., Emezue, C., Gangal, V., Garbacea, C., Hashimoto, T., Hou, Y., Jernite, Y., … Zhou, J. (2021). The GEM benchmark: Natural language generation, its evaluation and metrics. http://arxiv.org/abs/2102.01672

Howcroft, D. M., Belz, A., Clinciu, M.-A., Gkatzia, D., Hasan, S. A., Mahamood, S., Mille, S., Miltenburg, E. van, Santhanam, S., & Rieser, V. (2020). Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. Proceedings of the 13th International Conference on Natural Language Generation, 169–182. https://www.aclweb.org/anthology/2020.inlg-1.23

Howcroft, D. M., & Rieser, V. (2021). What happens if you treat ordinal ratings as interval data? Human evaluations in NLP are even more under-powered than you think. Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, 8932–8939. https://doi.org/10.18653/v1/2021.emnlp-main.703

Kamp, H., & Reyle, U. (2013). From discourse to logic: Introduction to modeltheoretic semantics of natural language, formal logic and discourse representation theory (Vol. 42). Springer Science & Business Media.

Mitchell, M., Wu, S., Zaldivar, A., Barnes, P., Vasserman, L., Hutchinson, B., Spitzer, E., Raji, I. D., & Gebru, T. (2019). Model cards for model reporting. Proceedings of the Conference on Fairness, Accountability, and Transparency, 220–229. https://doi.org/10.1145/3287560.3287596

Shimorina, A., & Belz, A. (2022). The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. Proceedings of the 2nd Workshop on Human Evaluation of Nlp Systems (Humeval), 54–75. https://aclanthology.org/2022.humeval-1.6

van der Lee, C., Gatt, A., Miltenburg, E. van, Wubben, S., & Krahmer, E. (2019). Best practices for the human evaluation of automatically generated text. Proceedings of the 12th International Conference on Natural Language Generation, 355–368. https://www.aclweb.org/anthology/W19-8643.pdf

van der Lee, C., Gatt, A., van Miltenburg, E., & Krahmer, E. (2021). Human evaluation of automatically generated text: Current trends and best practice

171

guidelines. Computer Speech & Language, 67, 101151.
https://doi.org/10.1016/j.csl.2020.101151

---

**Section 1:**  Paper and supplementary resources

Sections 1.1–1.3 record bibliographic and related information. These are
straightforward and don't warrant much in-depth explanation.

---

**Section 1.1:**  Details of paper reporting the evaluation experiment

---

**Question 1.1.1:**  Link to paper reporting the evaluation experiment.
Enter a link to an online copy of the the main reference (e.g., a paper) for the human
evaluation experiment. If the experiment hasn't been run yet, and the form is being
completed for the purpose of submitting it for preregistration, simply enter 'for
preregistration'.

> https://aclanthology.org/2022.acl-long.271.pdf

**Question 1.1.2:**  Which experiment within the paper is this form being
completed for?
Enter details of the experiment within the paper for which this sheet is being
completed. For example, the title of the experiment and/or a section number. If there is
only one human human evaluation, still enter the same information. If this is form is
being completed for pre-registration, enter a note that differetiates this experiment
from any others that you are carrying out as part of the same overall work.

> Human evaluation i.e. manual error annotation of redundancy for six
> data-to-text systems (described in section 7.2).

172

**Section 1.2:** Link to resources

---

**Question 1.2.1:** Link(s) to website(s) providing resources used in the evaluation experiment.

Enter the link(s). Such resources include system outputs, evaluation tools, etc. If there aren't any publicly shared resources (yet), enter 'N/A'.

> https://github.com/kasnerz/zeroshot-d2t-pipeline/
> (Only partial, full annotations provided via email.)

---

**Section 1.3:** Contact details

This section records the name, affiliation, and email address of person completing this sheet, and of the contact author if different.

---

**Section 1.3.1:** Details of the person completing this sheet.

---

**Question 1.3.1.1:** Name of the person completing this sheet.

Enter the name of the person completing this sheet.

> Filip Klubička

**Question 1.3.1.2:** Affiliation of the person completing this sheet.

Enter the affiliation of the person completing this sheet.

> ADAPT Centre, Technological University Dublin

**Question 1.3.1.3:** Email address of the person completing this sheet.

173

Enter the email address of the person completing this sheet.

> filip.klubicka@tudublin.ie

---

**Section 1.3.2:** Details of the contact author

---

**Question 1.3.2.1:** Name of the contact author.

Enter the name of the contact author, enter N/A if it is the same person as in Question 1.3.1.1

> N/A

---

**Question 1.3.2.2:** Affiliation of the contact author.

Enter the affiliation of the contact author, enter N/A if it is the same person as in Question 1.3.1.2

> N/A

---

**Question 1.3.2.3:** Email address of the contact author.

Enter the email address of the contact author, enter N/A if it is the same person as in Question 1.3.1.3

> N/A

---

**Section 2:** System Questions

Questions 2.1–2.5 record information about the system(s) (or human-authored stand-ins) whose outputs are evaluated in the Evaluation experiment that this sheet is being completed for. The input, output, and task questions in this section are closely interrelated: the value for one partially determines the others,as indicated for some combinations in Question 2.3.

174

**Question 2.1:** What type of input do the evaluated system(s) take?

This question is about the type(s) of input, where input refers to the representations and/or data structures shared by all evaluated systems. This question is about input type, regardless of number. E.g. if the input is a set of documents, you would still select text: document below.

Select all that apply. If none match, select 'other' and describe.

- ☑ 1. raw/structured data ⓘ
- ☐ 2. deep linguistic representation (DLR) ⓘ
- ☐ 3. shallow linguistic representation (SLR) ⓘ
- ☐ 4. text: subsentential unit of text ⓘ
- ☐ 5. text: sentence ⓘ
- ☐ 6. text: multiple sentences ⓘ
- ☐ 7. text: document ⓘ
- ☐ 8. text: dialogue ⓘ
- ☐ 9. text: other (please describe) ⓘ
- ☐ 10. speech ⓘ
- ☐ 11. visual ⓘ
- ☐ 12. multi-modal ⓘ
- ☐ 13. control feature ⓘ
- ☐ 14. no input (human generation) ⓘ
- ☐ 15. other (please describe) ⓘ

**Question 2.2:** What type of output do the evaluated system(s) generate?

This question is about the type(s) of output, where output refers to the and/or data structures shared by all evaluated systems. This question is about output type, regardless of number. E.g. if the output is a set of documents, you would still select *text: document* below. Note that the options for outputs are the same as for inputs except that the *no input (human generation) option* is replaced with *human-generated 'outputs'*, and the *control feature* option is removed.

Select all that apply. If none match, select 'other' and describe.

- ☐ 1. raw/structured data ⓘ
- ☐ 2. deep linguistic representation (DLR) ⓘ
- ☐ 3. Shallow linguistic representation (SLR) ⓘ
- ☐ 4. text: subsentential unit of text ⓘ
- ☑ 5. text: sentence ⓘ

175

☑ 6. text: multiple sentences  ⓘ

☐ 7. text: document  ⓘ

☐ 8. text: dialogue  ⓘ

☐ 9. text: other (please describe)  ⓘ

☐ 10. speech  ⓘ

☐ 11. visual  ⓘ

☐ 12. multi-modal  ⓘ

☐ 13. human generated 'outputs'  ⓘ

☐ 14. other (please describe)  ⓘ

---

**Question 2.3:** How would you describe the task that the evaluated system(s) perform in mapping the inputs in Q2.1 to the outputs in Q2.2?

This question is about the task(s) performed by the system(s) being evaluated. This is independent of the application domain (financial reporting, weather forecasting, etc.), or the specific method (rule-based, neural, etc.) implemented in the system. We indicate mutual constraints between inputs, outputs and task for some of the options below.

Occasionally, more than one of the options below may apply. Select all that apply. If none match, select 'other' and describe.

☐ 1. content selection/determination  ⓘ

☐ 2. content ordering/structuring  ⓘ

☐ 3. aggregation  ⓘ

☐ 4. referring expression generation  ⓘ

☐ 5. lexicalisation  ⓘ

☐ 6. deep generation  ⓘ

☐ 7. surface realisation (SLR to text)  ⓘ

☐ 8. feature-controlled text generation  ⓘ

☑ 9. data-to-text generation  ⓘ

☐ 10. dialogue turn generation  ⓘ

☐ 11. question generation  ⓘ

☐ 12. question answering  ⓘ

☐ 13. paraphrasing/lossless simplification  ⓘ

☐ 14. compression/lossy simplification  ⓘ

☐ 15. machine translation  ⓘ

☐ 16. summarisation (text-to-text)  ⓘ

☐ 17. end-to-end text generation  ⓘ

176

- [ ] 18. image/video description  ⓘ
- [ ] 19. post-editing/correction  ⓘ
- [ ] 20. other (please describe)  ⓘ

---

**Question 2.4:**  What are the input languages that are used by the system?

This question is about the language(s) of the inputs accepted by the system(s) being evaluated. Select any language name(s) that apply, mapped to standardised full language names in <u>ISO 639-1</u> (2019). E.g. English, Herero, Hindi. If no language is accepted as (part of) the input, select 'N/A'.

Select all that apply. If any languages you are using are not covered by this list, select 'other' and describe.

- [ ] 1. Abkhazian  ⓘ
- [ ] 2. Afar
- [ ] 3. Afrikaans
- [ ] 4. Akan
- [ ] 5. Albanian
- [ ] 6. Amharic
- [ ] 7. Arabic
- [ ] 8. Aragonese
- [ ] 9. Armenian
- [ ] 10. Assamese
- [ ] 11. Avaric  ⓘ
- [ ] 12. Avestan  ⓘ
- [ ] 13. Aymara
- [ ] 14. Azerbaijani  ⓘ
- [ ] 15. Bambara
- [ ] 16. Bashkir
- [ ] 17. Basque
- [ ] 18. Belarusian
- [ ] 19. Bengali  ⓘ
- [ ] 20. Bislama  ⓘ
- [ ] 21. Bosnian
- [ ] 22. Breton
- [ ] 23. Bulgarian
- [ ] 24. Burmese  ⓘ

177

- [ ] 25. Catalan, Valencian
- [ ] 26. Chamorro
- [ ] 27. Chechen
- [ ] 28. Chichewa, Chewa, Nyanja
- [ ] 29. Chinese
- [ ] 30. Church Slavic, Old Slavonic, Church Slavonic, Old Bulgarian, Old Church Slavonic ⓘ
- [ ] 31. Chuvash
- [ ] 32. Cornish
- [ ] 33. Corsican
- [ ] 34. Cree
- [ ] 35. Croatian
- [ ] 36. Czech
- [ ] 37. Danish
- [ ] 38. Divehi, Dhivehi, Maldivian
- [ ] 39. Dutch, Flemish ⓘ
- [ ] 40. Dzongkha
- [x] 41. English
- [ ] 42. Esperanto ⓘ
- [ ] 43. Estonian
- [ ] 44. Ewe
- [ ] 45. Faroese
- [ ] 46. Fijian
- [ ] 47. Finnish
- [ ] 48. French
- [ ] 49. Western Frisian ⓘ
- [ ] 50. Fulah ⓘ
- [ ] 51. Gaelic, Scottish Gaelic
- [ ] 52. Galician
- [ ] 53. Ganda
- [ ] 54. Georgian
- [ ] 55. German
- [ ] 56. Greek, Modern (1453–)
- [ ] 57. Kalaallisut, Greenlandic
- [ ] 58. Guarani
- [ ] 59. Gujarati

178

- [ ] 60. Haitian, Haitian Creole
- [ ] 61. Hausa
- [ ] 62. Hebrew ⓘ
- [ ] 63. Herero
- [ ] 64. Hindi
- [ ] 65. Hiri Motu
- [ ] 66. Hungarian
- [ ] 67. Icelandic
- [ ] 68. Ido ⓘ
- [ ] 69. Igbo
- [ ] 70. Indonesian
- [ ] 71. Interlingua (International Auxiliary Language Association) ⓘ
- [ ] 72. Interlingue, Occidental ⓘ
- [ ] 73. Inuktitut
- [ ] 74. Inupiaq
- [ ] 75. Irish
- [ ] 76. Italian
- [ ] 77. Japanese
- [ ] 78. Javanese
- [ ] 79. Kannada
- [ ] 80. Kanuri
- [ ] 81. Kashmiri
- [ ] 82. Kazakh
- [ ] 83. Central Khmer ⓘ
- [ ] 84. Kikuyu, Gikuyu
- [ ] 85. Kinyarwanda
- [ ] 86. Kirghiz, Kyrgyz
- [ ] 87. Komi
- [ ] 88. Kongo
- [ ] 89. Korean
- [ ] 90. Kuanyama, Kwanyama
- [ ] 91. Kurdish
- [ ] 92. Lao
- [ ] 93. Latin ⓘ
- [ ] 94. Latvian
- [ ] 95. Limburgan, Limburger, Limburgish

179

- [ ] 96. Lingala
- [ ] 97. Lithuanian
- [ ] 98. Luba-Katanga ⓘ
- [ ] 99. Luxembourgish, Letzeburgesch
- [ ] 100. Macedonian
- [ ] 101. Malagasy
- [ ] 102. Malay
- [ ] 103. Malayalam
- [ ] 104. Maltese
- [ ] 105. Manx
- [ ] 106. Maori ⓘ
- [ ] 107. Marathi ⓘ
- [ ] 108. Marshallese
- [ ] 109. Mongolian
- [ ] 110. Nauru ⓘ
- [ ] 111. Navajo, Navaho
- [ ] 112. North Ndebele ⓘ
- [ ] 113. South Ndebele ⓘ
- [ ] 114. Ndonga
- [ ] 115. Nepali
- [ ] 116. Norwegian
- [ ] 117. Norwegian Bokmål
- [ ] 118. Norwegian Nynorsk
- [ ] 119. Sichuan Yi, Nuosu ⓘ
- [ ] 120. Occitan
- [ ] 121. Ojibwa ⓘ
- [ ] 122. Oriya ⓘ
- [ ] 123. Oromo
- [ ] 124. Ossetian, Ossetic
- [ ] 125. Pali ⓘ
- [ ] 126. Pashto, Pushto
- [ ] 127. Persian ⓘ
- [ ] 128. Polish
- [ ] 129. Portuguese
- [ ] 130. Punjabi, Panjabi

180

- [ ] 131. Quechua
- [ ] 132. Romanian, Moldavian, Moldovan
- [ ] 133. Romansh
- [ ] 134. Rundi ⓘ
- [ ] 135. Russian
- [ ] 136. Northern Sami
- [ ] 137. Samoan
- [ ] 138. Sango
- [ ] 139. Sanskrit ⓘ
- [ ] 140. Sardinian
- [ ] 141. Serbian
- [ ] 142. Shona
- [ ] 143. Sindhi
- [ ] 144. Sinhala, Sinhalese
- [ ] 145. Slovak
- [ ] 146. Slovenian ⓘ
- [ ] 147. Somali
- [ ] 148. Southern Sotho
- [ ] 149. Spanish, Castilian
- [ ] 150. Sundanese
- [ ] 151. Swahili
- [ ] 152. Swati ⓘ
- [ ] 153. Swedish
- [ ] 154. Tagalog
- [ ] 155. Tahitian ⓘ
- [ ] 156. Tajik
- [ ] 157. Tamil
- [ ] 158. Tatar
- [ ] 159. Telugu
- [ ] 160. Thai
- [ ] 161. Tibetan ⓘ
- [ ] 162. Tigrinya
- [ ] 163. Tonga (Tonga Islands) ⓘ
- [ ] 164. Tsonga
- [ ] 165. Tswana
- [ ] 166. Turkish

181

- [ ] 167. Turkmen
- [ ] 168. Twi
- [ ] 169. Uighur, Uyghur
- [ ] 170. Ukrainian
- [ ] 171. Urdu
- [ ] 172. Uzbek
- [ ] 173. Venda
- [ ] 174. Vietnamese
- [ ] 175. Volapük ⓘ
- [ ] 176. Walloon
- [ ] 177. Welsh
- [ ] 178. Wolof
- [ ] 179. Xhosa
- [ ] 180. Yiddish
- [ ] 181. Yoruba
- [ ] 182. Zhuang, Chuang
- [ ] 183. Zulu
- [ ] 184. Other (please describe) ⓘ
- [ ] 185. N/A (please describe) ⓘ

---

**Question 2.5:** What are the output languages that are used by the system?

This field question the language(s) of the outputs generated by the system(s) being evaluated. Select any language name(s) that apply, mapped to standardised full language names in [ISO 639-1](#) (2019). E.g. English, Herero, Hindi. If no language is generated, select 'N/A'.

Select all that apply. If any languages you are using are not covered by this list, select 'other' and describe.

- [ ] 1. Abkhazian ⓘ
- [ ] 2. Afar
- [ ] 3. Afrikaans
- [ ] 4. Akan
- [ ] 5. Albanian
- [ ] 6. Amharic
- [ ] 7. Arabic
- [ ] 8. Aragonese
- [ ] 9. Armenian

- [ ] 10. Assamese
- [ ] 11. Avaric  *(i)*
- [ ] 12. Avestan  *(i)*
- [ ] 13. Aymara
- [ ] 14. Azerbaijani  *(i)*
- [ ] 15. Bambara
- [ ] 16. Bashkir
- [ ] 17. Basque
- [ ] 18. Belarusian
- [ ] 19. Bengali  *(i)*
- [ ] 20. Bislama  *(i)*
- [ ] 21. Bosnian
- [ ] 22. Breton
- [ ] 23. Bulgarian
- [ ] 24. Burmese  *(i)*
- [ ] 25. Catalan, Valencian
- [ ] 26. Chamorro
- [ ] 27. Chechen
- [ ] 28. Chichewa, Chewa, Nyanja
- [ ] 29. Chinese
- [ ] 30. Church Slavic, Old Slavonic, Church Slavonic, Old Bulgarian, Old Church Slavonic  *(i)*
- [ ] 31. Chuvash
- [ ] 32. Cornish
- [ ] 33. Corsican
- [ ] 34. Cree
- [ ] 35. Croatian
- [ ] 36. Czech
- [ ] 37. Danish
- [ ] 38. Divehi, Dhivehi, Maldivian
- [ ] 39. Dutch, Flemish  *(i)*
- [ ] 40. Dzongkha
- [x] 41. English
- [ ] 42. Esperanto  *(i)*
- [ ] 43. Estonian
- [ ] 44. Ewe

183

- [ ] 45. Faroese
- [ ] 46. Fijian
- [ ] 47. Finnish
- [ ] 48. French
- [ ] 49. Western Frisian ⓘ
- [ ] 50. Fulah ⓘ
- [ ] 51. Gaelic, Scottish Gaelic
- [ ] 52. Galician
- [ ] 53. Ganda
- [ ] 54. Georgian
- [ ] 55. German
- [ ] 56. Greek, Modern (1453–)
- [ ] 57. Kalaallisut, Greenlandic
- [ ] 58. Guarani
- [ ] 59. Gujarati
- [ ] 60. Haitian, Haitian Creole
- [ ] 61. Hausa
- [ ] 62. Hebrew ⓘ
- [ ] 63. Herero
- [ ] 64. Hindi
- [ ] 65. Hiri Motu
- [ ] 66. Hungarian
- [ ] 67. Icelandic
- [ ] 68. Ido ⓘ
- [ ] 69. Igbo
- [ ] 70. Indonesian
- [ ] 71. Interlingua (International Auxiliary Language Association) ⓘ
- [ ] 72. Interlingue, Occidental ⓘ
- [ ] 73. Inuktitut
- [ ] 74. Inupiaq
- [ ] 75. Irish
- [ ] 76. Italian
- [ ] 77. Japanese
- [ ] 78. Javanese
- [ ] 79. Kannada

184

- [ ] 80. Kanuri
- [ ] 81. Kashmiri
- [ ] 82. Kazakh
- [ ] 83. Central Khmer  ⓘ
- [ ] 84. Kikuyu, Gikuyu
- [ ] 85. Kinyarwanda
- [ ] 86. Kirghiz, Kyrgyz
- [ ] 87. Komi
- [ ] 88. Kongo
- [ ] 89. Korean
- [ ] 90. Kuanyama, Kwanyama
- [ ] 91. Kurdish
- [ ] 92. Lao
- [ ] 93. Latin  ⓘ
- [ ] 94. Latvian
- [ ] 95. Limburgan, Limburger, Limburgish
- [ ] 96. Lingala
- [ ] 97. Lithuanian
- [ ] 98. Luba-Katanga  ⓘ
- [ ] 99. Luxembourgish, Letzeburgesch
- [ ] 100. Macedonian
- [ ] 101. Malagasy
- [ ] 102. Malay
- [ ] 103. Malayalam
- [ ] 104. Maltese
- [ ] 105. Manx
- [ ] 106. Maori  ⓘ
- [ ] 107. Marathi  ⓘ
- [ ] 108. Marshallese
- [ ] 109. Mongolian
- [ ] 110. Nauru  ⓘ
- [ ] 111. Navajo, Navaho
- [ ] 112. North Ndebele  ⓘ
- [ ] 113. South Ndebele  ⓘ
- [ ] 114. Ndonga
- [ ] 115. Nepali

185

- 116. Norwegian
- 117. Norwegian Bokmål
- 118. Norwegian Nynorsk
- 119. Sichuan Yi, Nuosu (i)
- 120. Occitan
- 121. Ojibwa (i)
- 122. Oriya (i)
- 123. Oromo
- 124. Ossetian, Ossetic
- 125. Pali (i)
- 126. Pashto, Pushto
- 127. Persian (i)
- 128. Polish
- 129. Portuguese
- 130. Punjabi, Panjabi
- 131. Quechua
- 132. Romanian, Moldavian, Moldovan
- 133. Romansh
- 134. Rundi (i)
- 135. Russian
- 136. Northern Sami
- 137. Samoan
- 138. Sango
- 139. Sanskrit (i)
- 140. Sardinian
- 141. Serbian
- 142. Shona
- 143. Sindhi
- 144. Sinhala, Sinhalese
- 145. Slovak
- 146. Slovenian (i)
- 147. Somali
- 148. Southern Sotho
- 149. Spanish, Castilian
- 150. Sundanese

186

- [ ] 151. Swahili
- [ ] 152. Swati  ⓘ
- [ ] 153. Swedish
- [ ] 154. Tagalog
- [ ] 155. Tahitian  ⓘ
- [ ] 156. Tajik
- [ ] 157. Tamil
- [ ] 158. Tatar
- [ ] 159. Telugu
- [ ] 160. Thai
- [ ] 161. Tibetan  ⓘ
- [ ] 162. Tigrinya
- [ ] 163. Tonga (Tonga Islands)  ⓘ
- [ ] 164. Tsonga
- [ ] 165. Tswana
- [ ] 166. Turkish
- [ ] 167. Turkmen
- [ ] 168. Twi
- [ ] 169. Uighur, Uyghur
- [ ] 170. Ukrainian
- [ ] 171. Urdu
- [ ] 172. Uzbek
- [ ] 173. Venda
- [ ] 174. Vietnamese
- [ ] 175. Volapük  ⓘ
- [ ] 176. Walloon
- [ ] 177. Welsh
- [ ] 178. Wolof
- [ ] 179. Xhosa
- [ ] 180. Yiddish
- [ ] 181. Yoruba
- [ ] 182. Zhuang, Chuang
- [ ] 183. Zulu
- [ ] 184. Other (please describe)  ⓘ
- [ ] 185. N/A (please describe)  ⓘ

187

**Section 3:** Sample of system outputs, evaluators, and experimental design

---

**Section 3.1:** Sample of system outputs

Questions 3.1.1–3.1.3 record information about the size of the sample of outputs (or human-authored stand-ins) evaluated per system, how the sample was selected, and what its statistical power is.

---

**Question 3.1.1:** How many system outputs (or other evaluation items) are evaluated per system in the evaluation experiment?

Enter the number of system outputs (or other evaluation items) that are evaluated per system by at least one evaluator in the experiment. For most experiments this should be an integer, although if the number of outputs varies please provide further details here.

100

---

**Question 3.1.2:** How are system outputs (or other evaluation items) selected for inclusion in the evaluation experiment?

Select one option. If none match, select 'other' and describe:
- ◯ 1. by an automatic random process  ⓘ
- ◯ 2. by an automatic random process but using stratified sampling over given properties  ⓘ
- ⬤ 3. by manual, arbitrary selection  ⓘ
- ◯ 4. by manual selection aimed at achieving balance or variety relative to given properties  ⓘ
- ◯ 5. other (please describe)  ⓘ

---

188

### Section 3.1.3: Statistical power of the sample size.

---

**Question 3.1.3.1:** What method was used to determine the the statistical power of the sample size?

Enter the name of the method used.

> None provided

---

**Question 3.1.3.2:** What is the statistical power of the sample size?

Enter the numerical results of a statistical power calculation on the output sample.

> None provided

---

**Question 3.1.3.3:** Where can other researchers find details of the script used?

Enter a link to the script used (or another way of identifying the script). See, e.g., Card et al. (2020), Howcroft & Rieser (2021).

> None provided

---

### Section 3.2: Evaluators

Questions 3.2.1–3.2.5 record information about the evaluators participating in the experiment.

**Question 3.2.1:** How many evaluators are there in this experiment?

Enter the total number of evaluators participating in the experiment, as an integer.

2

**Section 3.2.2:** Evaluator Type

**Question 3.2.3:** How are evaluators recruited?

Please explain how your evaluators are recruited. Do you send emails to a given list? Do you post invitations on social media? Posters on university walls? Were there any gatekeepers involved? What are the exclusion/inclusion criteria?

Given the highly specific skillset required (PhD student or graduate-level NLP researcher) we reached out to reliable colleagues who we knew would be interested and would do a good job.

**Question 3.2.4:** What training and/or practice are evaluators given before starting on the evaluation itself?

Use this space to describe any training evaluators were given as part of the experiment to prepare them for the evaluation task, including any practice evaluations they did. This includes any introductory explanations they're given, e.g. on the start page of an online evaluation tool.

We sent them brief instructions via email and a definition of the redundancy quality criterion. Annotator training and guidelines were minimal to mirror the setting in the original study.

**Question 3.2.5:** What other characteristics do the evaluators have? Known either because these were qualifying criteria, or from information gathered as part of the evaluation.

Use this space to list any characteristics not covered in previous questions that the

190

evaluators are known to have, either because evaluators were selected on the basis of a characteristic, or because information about a characteristic was collected as part of the evaluation. This might include geographic location of IP address, educational level, or demographic information such as gender, age, etc. Where characteristics differ among evaluators (e.g. gender, age, location etc.), also give numbers for each subgroup.

> Key characteristic was their proficiency in English, their background in linguistics and NLO and their PhD-researcher-or-above academic level.

## Section 3.3: Experimental Design

Sections 3.3.1–3.3.8 record information about the experimental design of the evaluation experiment.

**Question 3.3.1:** Has the experimental design been preregistered? If yes, on which registry?

Select 'Yes' or 'No'; if 'Yes' also give the name of the registry and a link to the registration page for the experiment.

- ○ 1. yes
- ● 2. no

**Question 3.3.2:** How are responses collected?
Describe here the method used to collect responses, e.g. paper forms, Google forms, SurveyMonkey, Mechanical Turk, CrowdFlower, audio/video recording, etc.

> Google Sheets spreadsheet.

191

### Section 3.3.3: Quality assurance

Questions 3.3.3.1 and 3.3.3.2 record information about quality assurance.

---

**Question 3.3.3.1:** What quality assurance methods are used to ensure evaluators and/or their responses are suitable?

If any methods other than those listed were used, select 'other', and describe why below. If no methods were used, select *none of the above* and enter 'No Method'

Select all that apply:

- ☐ 1. evaluators are required to be native speakers of the language they evaluate. ⓘ
- ☐ 2. automatic quality checking methods are used during/post evaluation ⓘ
- ☐ 3. manual quality checking methods are used during/post evaluation ⓘ
- ☐ 4. evaluators are excluded if they fail quality checks (often or badly enough) ⓘ
- ☐ 5. some evaluations are excluded because of failed quality checks ⓘ
- ☐ 6. other (please describe) ⓘ
- ☑ 7. none of the above ⓘ

Please describe:

> The task was fairly rudimentary and required little quality assurance. There was a discussion step between the annotators after the annotation to agree on edge cases.

Please provide further details for your above selection(s)

---

**Question 3.3.3.2:** Please describe in detail the quality assurance methods that were used.
If no methods were used, enter 'N/A'

N/A

---

**Section 3.3.3:** Form/Interface

Questions 3.3.4.1 and 3.4.3.2 record information about the form or user interface that was shown to participants.

---

**Question 3.3.4.1:** Please include a link to online copies of the form/interface that was shown to participants.

Please record a link to a screenshot or copy of the form if possible. If there are many files, please create a signpost page (e.g., on GitHub that contains links to all applicable resouces). If there is a separate introductory interface/page, include it under Question 3.2.4.

https://docs.google.com/spreadsheets/d/15krRgujelUVWBLRn96

---

**Question 3.3.4.2:** What do evaluators see when carrying out evaluations?

Describe what evaluators are shown, in addition to providing the links in 3.3.4.1.

The sentence generated by the system and a field to note down the redundancy error counts.

---

**Question 3.3.5:** How free are evaluators regarding when and how quickly to carry out evaluations?

193

Select all that apply:

☑ 1. evaluators have to complete each individual assessment within a set time  ⓘ

☐ 2. evaluators have to complete the whole evaluation in one sitting  ⓘ

☐ 3. neither of the above (please describe)  ⓘ

**Question 3.3.6:** Are evaluators told they can ask questions about the evaluation and/or provide feedback?

Select all that apply.

☑ 1. evaluators are told they can ask any questions during/after receiving initial training/instructions, and before the start of the evaluation  ⓘ

☑ 2. evaluators are told they can ask any questions during the evaluation  ⓘ

☐ 3. evaluators are asked for feedback and/or comments after the evaluation, e.g. via an exit questionnaire or a comment box  ⓘ

☐ 4. other (please describe)  ⓘ

☐ 5. None of the above  ⓘ

**Question 3.3.7:** What are the experimental conditions in which evaluators carry out the evaluations?

Multiple-choice options (select one). If none match, select 'other' and describe.

🔘 1. evaluation carried out by evaluators at a place of their own choosing, e.g. online, using a paper form, etc.  ⓘ

◯ 2. evaluation carried out in a lab, and conditions are the same for each evaluator  ⓘ

◯ 3. evaluation carried out in a lab, and conditions vary for different evaluators  ⓘ

◯ 4. evaluation carried out in a real-life situation, and conditions are the same for each evaluator  ⓘ

◯ 5. evaluation carried out in a real-life situation, and conditions vary for different evaluators  ⓘ

194

○ 6. evaluation carried out outside of the lab, in a situation designed to resemble a real-life situation, and conditions are the same for each evaluator  ⓘ

○ 7. evaluation carried out outside of the lab, in a situation designed to resemble a real-life situation, and conditions vary for different evaluators  ⓘ

○ 8. other (please describe)  ⓘ

---

**Question 3.3.8:** Briefly describe the (range of different) conditions in which evaluators carry out the evaluations.

Use this space to describe the variations in the conditions in which evaluators carry out the evaluation, for both situations where those variations are controlled, and situations where they are not controlled. If the evaluation is carried out at a place of the evaluators' own choosing, enter 'N/A'

> On a laptop or computer, either at home or at university.

---

**Section 4:**  Quality Criteria – Definition and Operationalisation

Questions in this section collect information about each quality criterion assessed in the single human evaluation experiment that this sheet is being completed for.

---

**Many Criteria :**  Quality Criterion - Definition and Operationalisation
In this section you can create named subsections for each criterion that is being evaluated. The form is then duplicated for each criterion. To create a criterion type its name in the field and press the *New* button, it will then appear on tab that will allow you to toggle the active criterion. To delete the current criterion press the *Delete current* button.

> Redundancy (English)

New    Delete Current

Redundancy (English)

**Question 4.3.9:** How are raw responses from participants aggregated or otherwise processed to obtain reported scores for this quality criterion?

Normally a set of separate assessments is collected from evaluators and is converted to the results as reported. Describe here the method(s) used in the conversion(s). E.g. macro-averages or micro-averages are computed from numerical scores to provide summary, per-system results. If no such method was used, enter 'N/A'.

Counted and summed using google sheets formulae.

**Question 4.3.10:** Method(s) used for determining effect size and significance of findings for this quality criterion.

Enter a list of methods used for calculating the effect size and significance of any results, both as reported in the paper given in Question 1.1, for this quality criterion. If none calculated, state 'None'.

None

**Section 5:** Ethics

The questions in this section relate to ethical aspects of the evaluation. Information can be entered in the text box provided, and/or by linking to a source where complete information can be found.

**Question 5.1:** Has the evaluation experiment this sheet is being completed for, or the larger study it is part of, been approved by a research ethics committee? If yes, which research ethics committee?

Typically, research organisations, universities and other higher-education institutions require

some form ethical approval before experiments involving human participants, however innocuous, are permitted to proceed. Please provide here the name of the body that approved the experiment, or state 'No' if approval has not (yet) been obtained.

Yes, it is covered under general approval of the TU Dublin research ethics committee.

**Question 5.2:** Do any of the system outputs (or human-authored stand-ins) evaluated, or do any of the responses collected, in the experiment contain personal data (as defined in GDPR Art. 4, §1: https://gdpr.eu/article-4-definitions)? If yes, describe data and state how addressed.

State 'No' if no personal data as defined by GDPR was recorded or collected, otherwise explain how conformity with GDPR requirements such as privacy and security was ensured, e.g. by linking to the (successful) application for ethics approval from Question 5.1.

No.

**Question 5.3:** Do any of the system outputs (or human-authored stand-ins) evaluated, or do any of the responses collected, in the experiment contain special category information (as defined in GDPR Art. 9, §1: https://gdpr.eu/article-9-processing-special-categories-of-personal-data-prohibited)? If yes, describe data and state how addressed.

State 'No' if no special-category data as defined by GDPR was recorded or collected, otherwise explain how conformity with GDPR requirements relating to special-category data was ensured, e.g. by linking to the (successful) application for ethics approval from Question 5.1.

197

No.

---

**Question 5.4:** Have any impact assessments been carried out for the evaluation experiment, and/or any data collected/evaluated in connection with it? If yes, summarise approach(es) and outcomes.

Use this box to describe any *ex ante* or *ex post* impact assessments that have been carried out in relation to the evaluation experiment, such that the assessment plan and process, well as the outcomes, were captured in written form. Link to documents if possible. Types of impact assessment include data protection impact assessments, e.g. under GDPR. Environmental and social impact assessment frameworks are also available.

No.

---

**All Form Errors**

## List of all errors

refresh list of all errors

Press the button to refresh the list of all errors.

# ReproHum#0043:

## Human Evaluation

## Reproducing Language Model as an Annotator: Exploring Dialogue Summarization on AMI Dataset

### Vivian Fresen, Mei-Shin Wu-Urbanek, Steffen Eger

Adesso SE/Crif GmbH, Independent Researcher, Natural Language Learning Group (NLLG)
University of Mannheim
vivian.fresen@adesso.de, wumeishin@gmail.com, steffen.eger@uni-mannheim.de

### Abstract

This study, part of the ReproHum Belz and Thomson (2024) project, a collaborative effort among researchers to replicate and assess experiments published in the natural language processing (NLP) literature, replicates and evaluates "Language Model as an Annotator: Exploring DialoGPT for Dialogue Summarization" by Feng et al. (2021). Using DialoGPT, BART, and PGN models, we assess dialogue summarization's informativeness on a scale of 1 to 5. Surprisingly, our findings diverge from the original study, with different models producing the highest-rated summaries. This discrepancy suggests limitations in reproducing the original results and underscores the need for further investigation into dataset selection and model effectiveness.

**Keywords:** keyword1, keyword2, keyword3

## 1. Introduction

Reproducibility in natural language processing (NLP) is crucial for reliability, to ensure that independent researchers can arrive at the same conclusions by following the original report's documentation. In NLP, reproducilibity extends beyond model training parameters and may involve the entire evaluation process leading to reported results. While reproducibility has been studied in NLP e.g., for automatic metrics or models (Fokkens et al., 2013; Post, 2018; Chen et al., 2022), there is a scarcity of work addressing human evaluation.

Human evaluation is particularly important, however, as human annotations most often provide the ground-truth against which NLP models are compared.

The work reported in this paper forms part of the ReproHum[1] project, which focuses on enhancing the documentation of human evaluation properties and evaluating the consistency between results obtained in reproduction studies and those of the original research Belz et al. (2023); Belz and Thomson (2024).

Our focus paper is Feng et al. (2021). We followed the paper's guidelines to reproduce the automatic summarization outputs by using DialogGPT (Zhang et al., 2020). To do so, we leveraged four PhD students to assess the generated texts. Our goal was to assess whether we could reproduce the original results along specific selected dimensions.

Our report is structured as follows: Section 2 presents the original study design, providing an overview of the paper's content. In Section 3, we detail the reproduction of the NLP evaluation, outlining the specifics of the evaluation process to be replicated. Section 4 presents and discusses the results of the reproduced evaluation in comparison to the original paper. Finally, Section 5 offers concluding remarks and outlines avenues for future research.

## 2. Original Study Design

The original study, conducted by Feng et al. (2021), investigates enhancements to automatic text summarization. The study employs DialogGPT as an unsupervised annotator, focusing on three annotation aspects: keyword extraction, redundancy detection, and topic segmentation in dialogues.

Using DialogGPT, the authors annotate the SAMSum dataset (Gliwa et al., 2019) and the AMI dataset (Carletta et al., 2006), both containing dialogues and corresponding summaries. Pre-trained sequence-to-sequence BART (Lewis et al., 2020) and non-pretrained PGN (See et al., 2017) models are then used to generate summaries for the datasets annotated with keyword extraction ($D_{KE}$), redundancy detection ($D_{RD}$), topic segmentation ($D_{TS}$), and all three annotations combined ($D_{ALL}$) on both SAMSum and AMI datasets. The resulting summaries are assessed both automatically and manually.

BART($D_{KE}$) demonstrates superior performance in the SAMSum dataset to the baseline and

---

[1] https://reprohum.github.io/

PGN models, emphasizing the importance of keyword retention for concise dialogues. Conversely, PGN($D_{RD}$) exhibits significant improvements in the AMI dataset, highlighting the necessity of redundancy detection.

The study aims to investigate whether incorporating DialogGPT as a component in text summarization, specifically through keyword extraction, redundancy detection, and topic segmentation, enhances the efficacy and potential improvements in dialogue summarization. This is achieved by comparing its results against established models using BLEU and ROUGE metrics. The corresponding human evaluation process aimed to assess the informativeness, conciseness, and coverage of dialogue summaries. A total of 100 dialogues from SAMSum and 10 meetings from AMI, along with their corresponding generated summaries, were randomly sampled for evaluation. Four human evaluators were hired to rate each summary on a scale of 1 to 5 for each metric, with higher scores indicating better performance (Feng et al., 2021).

## 3. Reproduction Study Details

We aimed to replicate the original study as precisely as possible. We used a subset of AMI dataset consists of ten dialogues, which is the identical material in Feng et al. (2021)[2]. The AMI Meeting Corpus is a rich multi-modal dataset containing approximately 100 hours of meeting recordings. It comprises both scripted scenario-based meetings, simulating design team collaborations, and naturally occurring meetings across various domains. The dataset includes audio, video, and transcript data, making it suitable for research in speech recognition, natural language processing, and human-computer interaction (Carletta et al., 2006). On the other hand, the SAMSum Corpus is a dataset designed specifically for abstractive dialogue summarization. It consists of chat dialogues that have been manually annotated with abstractive summaries. The corpus serves as a benchmark for evaluating automated summarization models tailored to the unique challenges posed by dialogue data (Gliwa et al., 2019) . The SAMSum Corpus offers a high-quality resource for researchers to develop and refine techniques for generating concise and informative summaries from conversational exchanges.

In our reproduction study, we focused solely on the AMI dataset and the informativeness criterion. By concentrating solely on one criterion, the reproduction experiment is simplified and easier to follow. Moreover, evaluating only one criterion

should enable human evaluators to better concentrate on the annotation task; including multiple dimensions might confuse the annotators and hinder their ability to distinguish between the various criteria. However, this approach may potentially lead to overlooking other important aspects of dialogue summarization, resulting in less comprehensive evaluation outcomes.

Additionally, we followed the authors' guidelines and annotation outputs to be evaluated using the original scripts, without altering the summaries for each model and corresponding dialogue.

### 3.1. Evaluators

For the human evaluation, we engaged four annotators: native Chinese PhD students with high proficiency in English, as in the original study. One annotator is a PhD student in NLP, while the other three are from the fields of Sociology and Social Change. Each annotator received generous compensation,[3] as stipulated in the guidelines, for approximately 1-2 hours of work. The time estimation was based on the scope of the annotation task, which adhered to ReproHum recommendations. We adhered to the standardized ReproHum procedure for determining fair remuneration.

### 3.2. Differences to original study

The original study does not specify which interface was used for the Human Evaluators. We distributed the annotations to the evaluators via Google Forms (see Fig. 2 in the appendix), following the requirements set by the ReproHum team. This ensured uniform conditions and consistent result outputs for all reproduction experiments during the final evaluations and analyses. However, the outputs of the models, along with error annotations, remained consistent with those used in the original experiment.

In the original study, 100 dialogues from SAMSum and 10 meetings from AMI, along with their respective generated summaries, were assessed for informativeness, conciseness, and coverage by each model.

Fig. 1 shows an example of an AMI meeting with one of its summaries, followed by the option to rate the informativeness of the summary generated by the model. We were given the instruction to focus on the 10 AMI meetings only when reproducing the human evaluation, potentially to reduce annotation costs.

The instructions, originally provided in Chinese, were included with minor modifications by the authors. In Appendix A, we list them in the way we

---

[2]The full AMI dataset is provided in the repository on GitHub https://github.com/xcfcode/PLM_annotator.

[3]We paid each of them 50 EUR as a flat compensation in the form of amazon giftcards.

presented them for the human evaluation experiment. This approach facilitated a streamlined setup and enhanced accessibility for the annotators and the final evaluation process.

This is to prevent any potential influence on the reported outcomes. All information and resources should be accessed from the common resources folder provided by the project team. Any inquiries should be directed only to the ReproHum project managers, who communicated with the authors of the work being reproduced on behalf of the project.

## 4. Reproduction Results

We present our human evaluation result together with the scores provided in Feng et al. (2021) in Table 1. The comprehensive human evaluation results of the original article can be found in Table 4 in Appendix B.

The methodology for calculating the informativeness scores in the final evaluation results was not specified in Feng et al. (2021). Therefore, we utilized Python and R to calculate the informativeness scores over 10 AMI dialogs using three different methods: mean, median, and mode. Additionally, we adopted Feng et al. (2021)'s approach of using Fleiss' kappa score for evaluating inter-annotator agreement in our study. The scripts to calculate the mean, median, mode, as well as Fleiss' kappa scores are provided in our github repository [4].

|  | Model | Original | Mean | Median | Mode |
|---|---|---|---|---|---|
|  | Golden | 4.70 | 2.4 | 2.5 | 3 |
| AMI | PGN | 2.92 | 2.18 | 2.0 | 2 |
|  | HMNet | $3.52^{\dagger}$ | 2.2 | 2.0 | 2 |
|  | PGN($D_{KE}$) | 3.20 | 2.18 | 2.0 | 2 |
|  | PGN($D_{RD}$) | 3.15 | $3.0^{\dagger\dagger}$ | 3.0 | 3 |
|  | PGN($D_{TS}$) | 3.05 | 2.27 | 2.0 | 1 |
|  | PGN($D_{ALL}$) | $3.33^{\dagger\dagger}$ | $2.52^{\dagger}$ | 3.0 | 3 |

Table 1: Human evaluation results from Feng et al. (2021) is provided in the 'Original'column. The informativeness result in the reproduction experiment is provided in the 'Mean', 'Median' and 'Mode'columns. The corresponding Fleiss' kappa scores in the original paper are 0.48. The Fleiss' kappa score of our reproduction experiment is 0.069.

**Findings Comparison** The original results presented in the paper by Feng et al. (2021) indicate that their method, which combines DialogGPT as an annotator with BART and PGN as summarization generators, achieved the highest scores. Particularly, the combination of DialogGPT Redundancy reduction ($D_{RD}$) with both BART and PGN

resulted in better scores for conciseness (another dimension of annotation not considered by us). Additionally, when combined with DialogGPT Topic Segmentation ($D_{TS}$),the model performed better in coverage. However, HMNet, a Hierarchical Memory Network,[5] attained the best scores in informativeness and coverage for the AMI dataset.

There is a decisive gap between the scores of generated summaries and the scores of the gold summaries in the original study, indicating the increased difficulty of the AMI dataset (Feng et al., 2021). However, we did not observe such a significant difference between the score of Gold standard and the informativeness scores of the AMI dataset in our experiment. In Fig. 1 we can see an example of an AMI meeting provided for the human evaluation experiment, and its summary with the respectively rating options for informativeness. In Section C of the appendix, we give some examples where our raters disagree with the raters of the original study.

Our result in contrast with the original study is shown in Table 1.

The Informativeness measure applied to the gold outputs demonstrates a significant coefficient of variation(CV*) of 64.59%, indicating substantial variability relative to the mean value of 3.55. The unbiased sample standard deviation of 2.038 highlights considerable dispersion around the mean within the dataset. However, due to the small sample size of 2, the reliability of the standard deviation as a measure of dispersion may be limited.

Table 2 shows the coefficient of variation (CV*) with the corresponding mean values. The CV* metric is adapted for small sample sizes, making it suitable for use even with the limited pairs of results one may have (Belz, 2022).

Table 2: Coefficient of Variation (CV*) with Mean

| Sample | Mean | CV* |
|---|---|---|
| 1 | 3.55 | 64.59 |
| 2 | 3.22 | 18.58 |
| 3 | 3.47 | 15.52 |
| 4 | 3.10 | 3.22 |
| 5 | 2.19 | 1.14 |
| 6 | 2.59 | 31.79 |
| 7 | 2.40 | 10.41 |

Table 3: *

Note: CV* denotes the Coefficient of Variation.

---

[4]https://github.com/vivianCF/HumanEvaluation.git

[5]HMNet is a state-of-the-art model designed for abstractive dialogue summarization. It leverages memory modules and hierarchical attention mechanisms to capture dialogue nuances effectively. By storing relevant information and attending to different dialogue levels, HMNet generates coherent and informative summaries that faithfully represent the input dialogue.
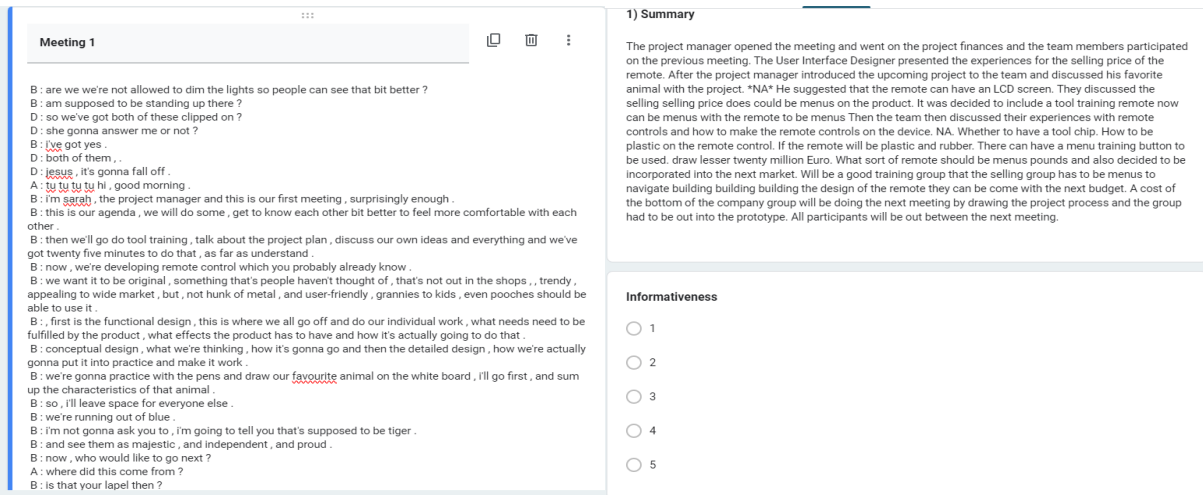
Figure 1: Example of an AMI meeting and its summary from one model.

Additionally, the wide confidence interval (-9.62 to 13.70) accentuates the uncertainty surrounding the true population mean, necessitating cautious interpretation of the dataset.

The coefficient of variation for PGN is 28.93%, indicating moderate variability relative to the mean value of 2.55. The sample standard deviation of 0.656 supports this observation, suggesting a moderate dispersion of data points around the mean. The confidence interval (-3.095, 4.407) implies some uncertainty about the true population mean. However, all measured values fall within one and two standard deviations from the mean, indicating a symmetric distribution around the mean.

HMNet exhibits a higher coefficient of variation at 46.02%, indicating high variability compared to PGN relative to the mean value of 2.86. The sample standard deviation of 1.170 suggests a greater dispersion of data points around the mean. The wider confidence interval (-5.521, 7.861) implies more uncertainty about the true population mean compared to Dataset 1. However, all measured values still fall within one and two standard deviations from the mean.

Similarly, $PGN(D_{KE})$ shows a coefficient of variation of 37.80%, indicating moderate variability compared to Dataset 1 relative to the mean value of 2.69. The sample standard deviation of 0.904 supports this, suggesting a moderate dispersion of data points around the mean. The confidence interval (-4.266, 6.074) also indicates some uncertainty about the true population mean. Nonetheless, like the other datasets, all measured values fall within one and two standard deviations from the mean.

$PGN(D_{RD})$ has the lowest coefficient of variation (4.86%), indicating the least variability compared to the mean value of 3.075 among all datasets. The sample standard deviation is also the small-est (0.133), suggesting minimal dispersion of data points around the mean.

In Contrast, $PGN(D_{TS})$ shows again a higher coefficient of variation (29.24%) compared to $PGN(D_{RD})$, indicating higher variability relative to the mean value of 2.66. The sample standard deviation is also larger (0.691), suggesting a greater dispersion of data points around the mean.

$PGN(D_{ALL})$ shows a coefficient of variation of 27.61%, slightly lower than $PGN(D_{TS})$, indicating slightly less variability compared to the mean value of 2.925. The sample standard deviation (0.718) is comparable to $PGN(D_{TS})$, suggesting a similar dispersion of data points around the mean.

In summary, $PGN(D_{RD})$ demonstrates the least variability, followed by $PGN(D_{ALL})$ and PGN, respectively. Overall, despite variations in coefficient of variation and sample size, all datasets exhibit symmetric distributions around the mean, as indicated by all measured values falling within one and two standard deviations from the mean.

Both the original study and our reproduction experiment suggest that PGN combined with redundancy reduction can achieve good performance for the AMI dataset in dialogue summarization. However, the gap between the gold standard and the other datasets in our case is not substantial, with a score of 2.4; the score is still considerably lower than the original result of 4.70.

In summary, these findings indicate a significant deviation compared to the original study. There are no significant discrepancies observed between the gold standards and the remaining datasets in our experiment, suggesting a different behavior compared to the original study. We observe a distinct trend compared to the original study; for instance, in our experiments, PGN $(D_{RD})$ demonstrates the highest performance, with PGN $(D_{ALL})$ closely fol-

lowing, which is somewhat unexpected considering that in the original study, HMNet achieved the highest score followed by PGN ($D_{ALL}$).

In our reproduction experiment, the scores were overall inferior to those in the original study, mostly all below 3.0 versus the original scores were all above 3.0. Specifically, the gold standard scores in our analysis are significantly lower than those reported in the original study. In our experiment, we have noticed a distinct trend that contrasts with the findings of Feng et al. (2021) in which the performance of HMNet does not exhibit substantial gains over the PGN models.

Upon comparing the coefficient of variation (CV*), it becomes evident that more replications of the same experiment may be required to draw more robust conclusions about the results presented in the human evaluation reproduction approach.

Furthermore, in our annotation task, we did not achieve comparable agreement (0.069), for AMI and informativeness on the same dataset. This is in strong contrast to the original study, which reported agreements of above 0.40 throughout.

## 5. Conclusion

Our research focuses on the reproduction and evaluation of dialogue summarization models through human assessment. The collaboration with the ReproHum organizers and access to materials from the original authors greatly facilitated the successful execution of our reproduction experiment.

**Our key findings include:**

- In our reproduction study, the inter-annotator agreement was notably lower, registering at 0.069, compared to above 0.40 reported in the original study.

- We were unable to confirm the effectiveness of the proposed approach in terms of informativeness. While we observed a moderate positive Pearson correlation coefficient of 0.481 between the informativeness of the original study and our experiment, indicating a medium level of correlation, the Spearman correlation coefficient of approximately -0.058 between both experiments suggests a weak negative monotonic relationship. Overall, the correlations between the original human evaluation and our reproduction are weak. However, it must be kept in mind that we correlated vectors of very short size (length seven).

  Notable discrepancies in human evaluation outcomes persist, indicating potential variations in annotators, methodology, or dataset

selection for each dialogue summarization. We hypothesize that these differences could be attributed to two main factors. Firstly, the dataset is limited to only 10 meeting materials, which may lead to skewed average scores, favoring extreme values. Additionally, the involvement of only three evaluators may not provide a comprehensive assessment. Further experiments and reproductions are necessary to draw more conclusive findings from this study.

- The evaluated model performances in our reproduction study were inferior compared to the scores reported in the original study.

- The ratings in the "Original" column of Table 1 were not elucidated in the original study. From the context, we inferred the authors reported the average scores. However, in a small-scale study, one method to prevent outliers from impacting the mean is to utilize the median. Consequently, we were unsure whether it is indicative of a mean or a median.

- The human evaluation in our case is neither repeatable nor reproducible.

A potential explanation for these results is the persistent gap between the scores of generated summaries and those of gold summaries for the AMI dataset, indicating its inherent difficulty. The complexity and ambiguity of the dialogues posed a challenge during the experiment's preparation, making them difficult to follow and leading to divergent ratings among evaluators.

Moreover, the original study regarding AMI/Informativeness did not demonstrate effectiveness: the baseline HMNet performed the best. This raises the question of whether the selection of the AMI dataset is appropriate for the human evaluation reproduction and the verification of the performance of the models using DialoGPT to achieve better performance in dialogue summarization. Furthermore, conducting a comprehensive analysis of dataset characteristics and evaluation metrics could offer valuable insights into enhancing the appropriateness of the dataset selection for evaluating summarization models.

Our reproduction study raises an intriguing question about the identification and management of subjective practices that might have been employed in the original study. The lack of information on human participants' training depth and the undisclosed time investment in annotations during the original study contribute to uncertainties in interpreting the significant disparity in our human evaluation results.

## Acknowledgements

## 6. Bibliographical References

Anya Belz. 2022. A metrological perspective on reproducibility in nlp. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *The Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

BSI. 1973a. *Natural Fibre Twines*, 3rd edition. British Standards Institution, London. BS 2570.

BSI. 1973b. Natural fibre twines. BS 2570, British Standards Institution, London. 3rd. edn.

Jean Carletta, Simone Ashby, Sebastien Bourban, Mike Flynn, Mael Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa Kronenthal, Guillaume Lathoud, Mike Lincoln, Agnes Lisowska, Iain McCowan, Wilfried Post, Dennis Reidsma, and Pierre Wellner. 2006. The ami meeting corpus: A pre-announcement. In *Machine Learning for Multimodal Interaction*, pages 28–39, Berlin, Heidelberg. Springer Berlin Heidelberg.

A. Castor and L. E. Pollux. 1992. The use of user modelling to guide inference and learning. *Applied Intelligence*, 2(1):37–53.

Jiaao Chen and Diyi Yang. 2020. Multi-view sequence-to-sequence models with conversational structure for abstractive dialogue summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4106–4118, Online. Association for Computational Linguistics.

Yanran Chen, Jonas Belouadi, and Steffen Eger. 2022. Reproducibility issues for BERT-based evaluation metrics. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2965–2989, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

J.L. Chercheur. 1994. *Case-Based Reasoning*, 2nd edition. Morgan Kaufman Publishers, San Mateo, CA.

N. Chomsky. 1973. Conditions on transformations. In *A festschrift for Morris Halle*, New York. Holt, Rinehart & Winston.

Umberto Eco. 1990. *The Limits of Interpretation*. Indian University Press.

Xiachong Feng, Xiaocheng Feng, Libo Qin, Bing Qin, and Ting Liu. 2021. Language model as an annotator: Exploring DialoGPT for dialogue summarization. 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 1479–1491, Online. Association for Computational Linguistics.

Antske Fokkens, Marieke van Erp, Marten Postma, Ted Pedersen, Piek Vossen, and Nuno Freire. 2013. Offspring from reproduction problems: What replication failure teaches us. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1691–1701, Sofia, Bulgaria. Association for Computational Linguistics.

Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. SAMSum corpus: A human-annotated dialogue dataset for abstractive summarization. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*,

pages 70–79, Hong Kong, China. Association for Computational Linguistics.

Paul Gerhard Hoel. 1971a. *Elementary Statistics*, 3rd edition. Wiley series in probability and mathematical statistics. Wiley, New York, Chichester. ISBN 0 471 40300.

Paul Gerhard Hoel. 1971b. *Elementary Statistics*, 3rd edition, Wiley series in probability and mathematical statistics, pages 19–33. Wiley, New York, Chichester. ISBN 0 471 40300.

Otto Jespersen. 1922. *Language: Its Nature, Development, and Origin*. Allen and Unwin.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Manling Li, Lingyu Zhang, Heng Ji, and Richard J. Radke. 2019. Keep meeting summaries on topic: Abstractive multi-modal meeting summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2190–2196, Florence, Italy. Association for Computational Linguistics.

Ian Magnusson, Noah A. Smith, and Jesse Dodge. 2023. Reproducibility in NLP: What have we learned from the checklist? In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12789–12811, Toronto, Canada. Association for Computational Linguistics.

Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing order into text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 404–411, Barcelona, Spain. Association for Computational Linguistics.

Matt Post. 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. Association for Computational Linguistics.

Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.

Guokan Shang, Wensi Ding, Zekun Zhang, Antoine Tixier, Polykarpos Meladianos, Michalis Vazirgiannis, and Jean-Pierre Lorré. 2018. Unsupervised abstractive meeting summarization with multi-sentence compression and budgeted submodular maximization. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 664–674, Melbourne, Australia. Association for Computational Linguistics.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Charles Joseph Singer, E. J. Holmyard, and A. R. Hall, editors. 1954–58. *A history of technology*. Oxford University Press, London. 5 vol.

Jannik Strötgen and Michael Gertz. 2012. Temporal tagging on different domains: Challenges, strategies, and gold standards. In *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, pages 3746–3753, Istanbul, Turkey. European Language Resource Association (ELRA).

S. Superman, B. Batman, C. Catwoman, and S. Spiderman. 2000. *Superheroes experiences with books*, 20th edition. The Phantom Editors Associates, Gotham City.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. DIALOGPT : Large-scale generative pre-training for conversational response generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278, Online. Association for Computational Linguistics.

Chenguang Zhu, Ruochen Xu, Michael Zeng, and Xuedong Huang. 2020. A hierarchical network for abstractive meeting summarization with cross-domain pretraining. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 194–203, Online. Association for Computational Linguistics.
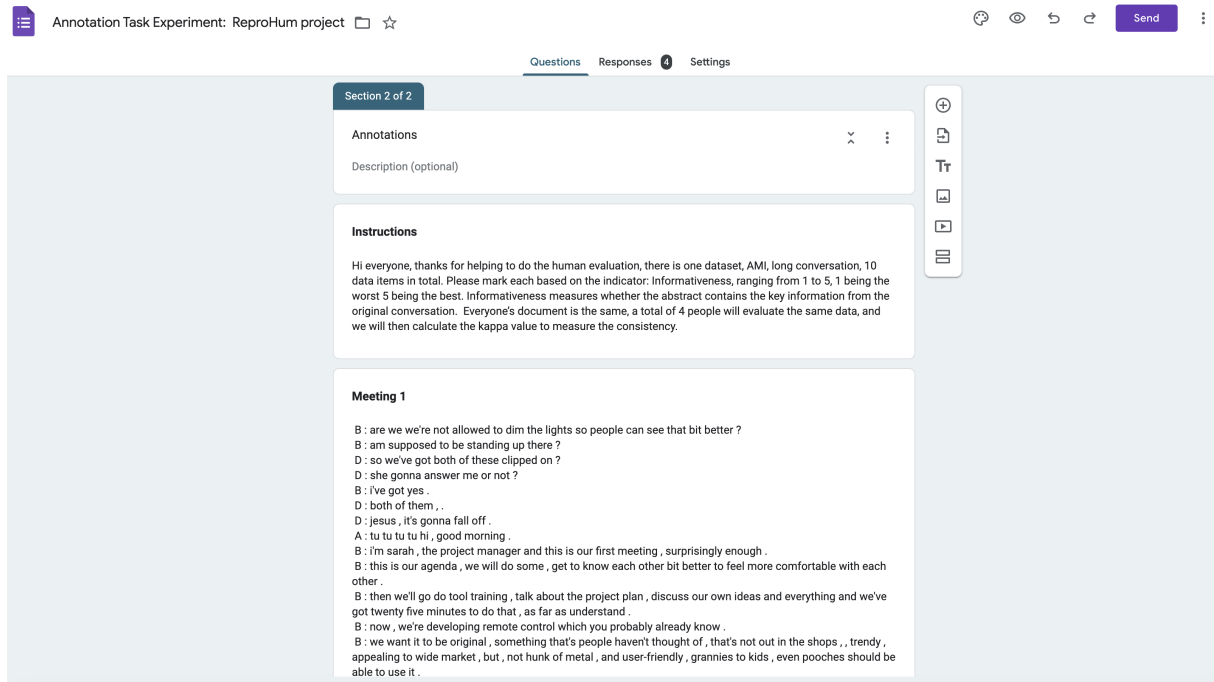
# Appendix

## A.  Annotator Guidelines



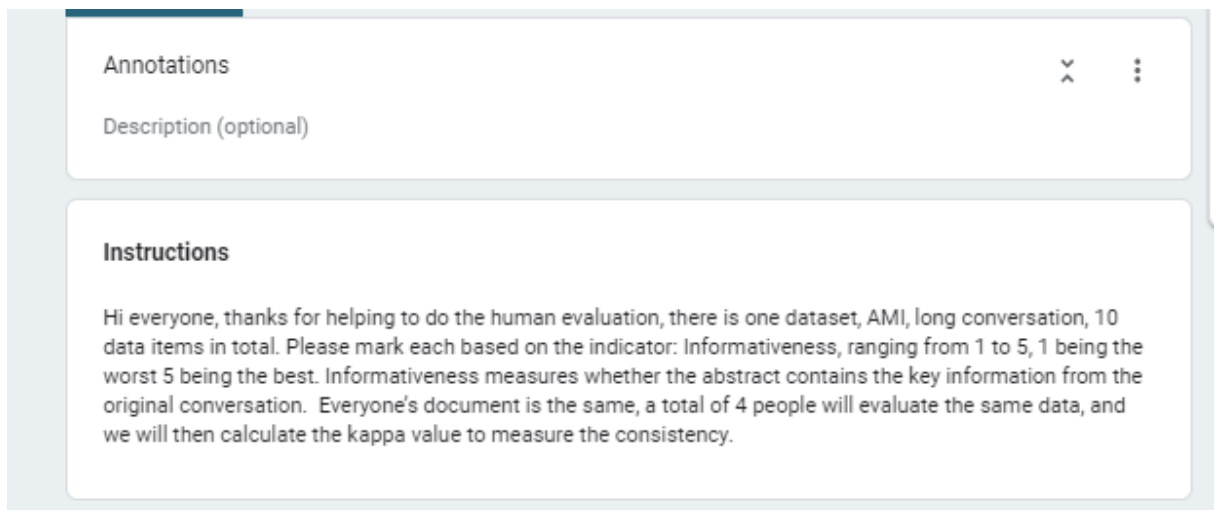Figure 2: Example of Google Forms interface used during the Human Evaluation reproduction experiment



Figure 3: Example of AMI instructions for Human Evaluation

## B. Results from original Study

| | Model | Info. | Conc. | Cov. |
|---|---|---|---|---|
| | Golden | **4.37** | 4.26 | 4.27 |
| SamSum | BART | 3.66 | 3.65 | 3.66 |
| | MV-BART | 3.85 | 3.76 | 3.88 |
| | BART($D_{KE}$) | 3.88 | 3.77 | 3.79 |
| | BART($D_{RD}$) | 3.74 | **3.98**$^{\dagger}$ | 3.89 |
| | BART($D_{TS}$) | **3.95**$^{\dagger}$ | 3.76 | **4.01**$^{\dagger\dagger}$ |
| | BART($D_{ALL}$) | **4.05**$^{\dagger}$ | **3.78**$^{\dagger\dagger}$ | **4.08**$^{\dagger}$ |
| | Golden | 4.70 | 3.85 | 4.35 |
| AMI | PGN | 2.92 | 3.08 | 2.70 |
| | HMNet | **3.52**$^{\dagger}$ | 2.40 | **3.40**$^{\dagger}$ |
| | PGN($D_{KE}$) | 3.20 | 3.08 | 3.00 |
| | PGN($D_{RD}$) | 3.15 | **3.25**$^{\dagger}$ | 3.00 |
| | PGN($D_{TS}$) | 3.05 | **3.10**$^{\dagger\dagger}$ | **3.17**$^{\dagger\dagger}$ |
| | PGN($D_{ALL}$) | **3.33**$^{\dagger\dagger}$ | **3.25**$^{\dagger}$ | 3.10 |

Table 4: Human evaluation results from the original paper indicate the following abbreviations: 'Info.' for informativeness, 'Conc.' for conciseness, and 'Cov.' for coverage. In the case of SAMSum, the inter-annotator agreement (Fleiss' kappa) scores for each metric are 0.46, 0.37, and 0.43, respectively. For AMI, the corresponding Fleiss' kappa scores are 0.48, 0.40, and 0.41.

| Fleiss' Kappa Value | Interpretation |
|---|---|
| 0.00 - 0.20 | Slight agreement |
| 0.21 - 0.40 | Fair agreement |
| 0.41 - 0.60 | Moderate agreement |
| 0.61 - 0.80 | Substantial agreement |
| 0.81 - 1.00 | Almost perfect agreement |

Table 5: Interpretation of Fleiss' Kappa Values

## C. Summary Examples



**2) Summary**

the project manager introduced the project to the team and then the team members participated in an exercise in which they drew their favorite animal on the white board and discussed why they liked the particular animal . the project manager discussed the project finances and selling prices . the group then evaluated the project process , and discussed their experiences with remote controls . they discussed making the remote universally compatible , to be a mobile phone , and using plastic instead of metal . the user interface designer discussed the controls and how they operated together , and presented the type of scroll in the shared design . the marketing expert led an evaluation of the prototype . the remote will be made of plastic and will feature chunky colors and voice . production costs cannot exceed 12 .50 euro . monkey control will be small , appealing to a wide market , but not a metal user youth . having a feature on the remote which allows the user to locate a lost seagull . having discussion about a dual function on the screen , channel up , two basic functions , and one for the basic functions of the user interface . she shares the list of the features that will be incorporated into the design of a remote that should be used . she will run out of the blue . the sheepdog will be a flying seagull . the marketing expert will work on trend watching the remote control will not contain a useful feature . will have a menu , display , menu , channel , and menu buttons . would like to see if it is misplaced by someone like it , so that the eagle was not a vampire bat and will have to be used with plastic . she had to use plastic for the display since the remote controls were too expensive . however , the group decided that they had to re it every time and to eliminate the signals . it might not be too costly to add to look good and incorporate a feature that has since the product is unlikely that the remote will make it desirable and be fashionable so that it can be a soft and stylish , fashionable and cool and cool , traditional remote .

Figure 4: Example of summary (HMNet) for the meeting 1 of AMI data set.

**3) Summary**

The project manager opened the meeting and introduced the upcoming project to the team members introduce themselves and then decided over the agenda. The team then discussed their experiences from the previous meeting. They discussed the project process and discussed the cost of the remote. *NA* *NA*. It will be made of rechargeable with a docking station and sleek theme. He discussed the extra buttons that should be incorporated into the design. After the project manager talked about the project finances and the team discussed the features they could be made from the next meeting and has them as them as they liked about the cost for the remote controls and they decided to incorporate a docking station. Then the team evaluated the project budget and the group had to focus with the project. Overall, the group decided to include a wide with the remote control. NA. Whether the user interface designer will look at the functional requirements and has a remote controls was hard buttons could be incorporated to the budget of the product. How to use a plastic which should be used. What sort of extra buttons or the remote will be made Will not have a wide station that should have a trendy station as they could have a teletext function. There will not have an LCD screen. Size of rechargeable screen and corporate image that they want to use up to the project budget. A extra buttons are more

Figure 5: Example of summary (PGN(D_KE)) for the meeting1 of AMI dataset.

**4) Summary**

The project manager opened the meeting by stating the agenda. The industrial designer discussed the interior workings of a remote which is easier to be He discussed the target group of the features of the function of the device. They also discussed using a timer be an LCD screen. After the project manager closes the meeting and going the team members introduce themselves by name and the team discussed their favorite animal and discussed what features they wanted to include for the remote they be easier to use. Then the team found the project finances and what features the target goals. *NA* NA. and at the functional functional group to be out from the working design. It was decided to include a timer to use a timer which can only be used for the television or the television remote which will only be set on the functional design, and the the marketing will be no limited at a regular target group to have a feature to address their product. There will have a hook screen which is no at address a feature of the remote. A decision that the remote would have a display function. Whether to have an LCD screen which will be used to address the and the television target group had to be no limited for the budget. When they have not sure their project process and the project manager's closing, the User Interface Designer and the Industrial Designer to research possible important and the group decided to use an LCD screen, and the Marketing Expert to prepare the working design of the remote and that they will include to include the remote which should be easier

Figure 6: Example of summary (PGN(D_TS)) for the meeting10 of AMI dataset.

**1) Summary**

The project manager opened the meeting and went on the project finances and the team members participated on the previous meeting. The User Interface Designer presented the experiences for the selling price of the remote. After the project manager introduced the upcoming project to the team and discussed his favorite animal with the project. *NA* He suggested that the remote can have an LCD screen. They discussed the selling selling price does could be menus on the product. It was decided to include a tool training remote now can be menus with the remote to be menus Then the team then discussed their experiences with remote controls and how to make the remote controls on the device. NA. Whether to have a tool chip. How to be plastic on the remote control. If the remote will be plastic and rubber. There can have a menu training button to be used. draw lesser twenty million Euro. What sort of remote should be menus pounds and also decided to be incorporated into the next market. Will be a good training group that the selling group has to be menus to navigate building building building the design of the remote they can be come with the next budget. A cost of the bottom of the company group will be doing the next meeting by drawing the project process and the group had to be out into the prototype. All participants will be out between the next meeting.

Figure 7: Example of summary (PGN(D_RD)) for the meeting10 of AMI dataset.

## D. HEDS

Below is the HEDS of the Human Evaluation Experiment. The original HEDS are based on the provided documents of the ReproHuman Group (Shimorina and Belz, 2022), which can be accessed at: `https://favorite-fox.static.domains/heds-2022-11-18`. For more information and updates, please visit the ReproNLP 2024 GitHub repository: `https://github.com/nlp-heds/repronlp2024`.

# ReproHum #0712-01: Human Evaluation Reproduction Report for "Hierarchical Sketch Induction for Paraphrase Generation"

## Mohammad Arvan and Natalie Parde

University of Illinois at Chicago

{marvan3,parde}@uic.edu

## Abstract

Human evaluations are indispensable in the development of NLP systems because they provide direct insights into how effectively these systems meet real-world needs and expectations. Ensuring the reproducibility of these evaluations is vital for maintaining credibility in natural language processing research. This paper presents our reproduction of the human evaluation experiments conducted by Hosking et al. (2022) for their paraphrase generation approach. Through careful replication we found that our results closely align with those in the original study, indicating a high degree of reproducibility.

**Keywords:** reproducibility, human evaluation, open science, paraphrase generation

## 1. Introduction

Human evaluation serves as the cornerstone for appraising the efficacy of machine learning and natural language processing pipelines. Consequently, understanding and addressing the challenges (Howcroft et al., 2020) that may impede the reproducibility of human evaluation experiments is paramount. The ReproHum Project (Belz and Thomson, 2024) is dedicated to devising a methodological framework specifically tailored to assess the reproducibility of human evaluation experiments within the domain of Natural Language Processing (NLP). In line with analogous meta-analytical endeavors (Open Science Collaboration, 2015; Errington et al., 2021a,b), this project seeks to heighten rigor, transparency, and reliability in NLP research. Furthermore, insights garnered from this initiative may help refine future human evaluation methodologies, enhancing their dependability and credibility.

ReproHum has been broken into multiple stages or *rounds*; we are presently near the end of round one. The primary objective of round one is to identify a set of experiments that are reproducible under the same conditions. Additional details regarding round one can be found in §2 and the process is also extensively reported by Belz et al. (2023). Work reported in this paper is part of the second batch of experiments selected for round one of ReproHum.

Specifically, we reproduced the human evaluation experiments conducted in the paper "Hierarchical Sketch Induction for Paraphrase Generation" by Hosking et al. (2022). The original study compared four models for paraphrase generation, and human evaluators assessed the quality of the generated paraphrases. Thanks to the cooperation of the original authors, we were able to reproduce the human evaluation experiments as closely as possible to

the original study. We compared our results to the original outcomes, finding that the results of our reproduction are very close to the originally reported results. This suggests that the human evaluation experiments conducted in the original study have a high degree of reproducibility. We have released the data, code, and results of our reproduction to ensure transparency and facilitate further research in this area.[1]

## 2. Background

In the first step of the ReproHum Project (Belz et al., 2023), 177 papers were identified that (a) contained human evaluation, and (b) were published in the *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)* or the *Transactions of the Association for Computational Linguistics (TACL)* in the 2018-2022 period. Through a multi-stage process, 20 experiments from 15 papers were selected to be reproduced. The selection process involved manual review for suitability, responsiveness of the original authors, and availability of a predetermined set of relevant details. Selected experiments were annotated with categorical labels indicating the number of evaluators (small, not small), cognitive complexity[2] (low, medium, high), and evaluator training and expertise (neither, either, both). Following annotation, six of the 20 selected experiments were chosen to achieve a balance of inclusion of these factors in the first batch of reproductions.

In **round one, batch a**, each included experiment was assigned to two partner labs. The partner labs were instructed to reproduce the experiment as closely as possible, given the information

---

[1] https://github.com/mo-arvan/paraphrase-generation-reproduction

[2] Based on scores given to each criterion in Appendix E of Howcroft et al. (2020).

provided in the original paper and any additional information and clarification obtained through direct communication between the original authors and the ReproHum leadership team. Partner labs were also instructed to document any deviations from the original experiment and the reasons for these deviations. The results of the reproduction were compared to the originally published results to assess the extent to which the experiment was reproducible. These reports and the corresponding data were published as part of the *2023 ReproNLP Shared Task on Reproducibility of Evaluations in NLP* (Belz and Thomson, 2023; González Corbelle et al., 2023; Watson and Gkatzia, 2023; Arvan and Parde, 2023; van Miltenburg et al., 2023; Ito et al., 2023; Gao et al., 2023; Mieskes and Benz, 2023; Hürlimann and Cieliebak, 2023; Platek et al., 2023; Klubička and Kelleher, 2023; Li et al., 2023; Mahamood, 2023).

Overall, the results from round one suggest a varied degree of reproducibility across the experiments, with some being easily reproduced and others not. By analyzing the attributes of each experiment and the corresponding results of the reproduction, it can be inferred that the higher the cognitive complexity, the lower the degree of reproducibility. The total number of evaluators also had an inverse correlation with the degree of reproducibility. While these preliminary findings are insightful, the ReproHum team acknowledged that they are based on a small sample size and may not be generalizable. Hence, another batch of experiments was selected for additional reproducibility assessment (**round one, batch b**). It is our round one, batch b results that we report in this paper.

# 3. Methods

For round one, batch b, we were assigned to reproduce "Hierarchical Sketch Induction for Paraphrase Generation" (Hosking et al., 2022). The ReproHum leadership team shared a document containing general instructions for reproduction and experiment-specific information for this paper. We summarize the paper and our methods for reproducing it below.

## 3.1. Hierarchical Sketch Induction for Paraphrase Generation

Hosking et al. (2022) introduced a new generative model called Hierarchical Refinement Quantized Variational Autoencoders (HRQ-VAE). Their proposed model utilized syntactic sketch for paraphrase generation, drawing parallels to the way humans plan out utterances and using those similarities in a sketching step added to the model to help in generating paraphrases. They evalu-

Informed Consent

This study is being conducted by researchers at the School of Informatics, University of Edinburgh. If you have any questions about this study, feel free to contact us (tom.hosking@ed.ac.uk). Participation in this research is voluntary. You have the right to withdraw from the experiment at any time. The collected data will be used for research purposes only. All output data will be anonymised and we will not collect or store any information that could be used to identify who you are. A full Participant Information Sheet is available here.

☐ I understand the participant information and consent to participate in this study.

If you do not consent, please return this HIT.

Instructions

In this task you will read roughly thirty examples of sentences and two paraphrases created by a computer program. The program aims to rewrite the sentence so that it means the same thing, but using different words and/or different word order.

**Please read all the sentences carefully**, this should take you about 20 minutes (if you do the task very quickly your HIT will be rejected).

You will be asked to choose which system performs better, for three aspects of the paraphrases:

1. Which system output is the most **fluent and grammatical**?
2. To what extent is the **meaning** expressed in the original sentence **preserved** in the rewritten version, with no additional information added?
3. Does the rewritten version use **different words or phrasing** to the original? You should choose the system that uses the **most** different words or word order.

Remember that you are being asked to **rate the system**, not the original.

Some of the sentences only have small differences! Be careful to choose the one that is **most different** for the dissimilarity category. If the control samples are not answered correctly then we will assume that you have answered at random and reject the HIT.

A small number of samples may have two choices that are *exactly* the same - in these cases please pick an answer at random, this will not cause the HIT to be rejected.

Examples

First, complete these example tasks correctly:



Figure 1: The user interface used for human evaluation in the original study.

ated the performance of their model compared to several baseline models on the Paralex (Fader et al., 2013), Quora Question Pairs (QQP)[3] and MSCOCO datasets (Lin et al., 2014).

For baselines, the authors compared their approach to Gaussian Variational AutoEncoder (VAE) (Bowman et al., 2016), Latent Bag-of-Words (Fu et al., 2019), Separator (Hosking and Lapata, 2021), and several other paraphrase generation systems. They evaluated their approach and the baselines on the mentioned datasets using iBLEU (Sun and Zhou, 2012), BLEU, Self-BLEU, and P-BLEU. iBLEU, the primary evaluation metric, is a variant of BLEU that uses a paraphrase dataset to evaluate paraphrase quality by assessing the faithfulness of generated outputs compared to ref-

---

[3]https://kaggle.com/competitions/quora-question-pairs

erence paraphrases. It also gauges the extent to which diversity is incorporated. Their automated evaluation suggests that the VAE, Latent BoW, Separator, and HRQ-VAE models performed best.

The four top-performing models were then selected for additional human evaluation. The human evaluation was conducted on Amazon Mechanical Turk (MTurk). It involved 180 human intelligence tasks (HITs), each containing 32 paraphrase pairs.[4] Each task contained two attention checks to ensure the quality of the responses. MTurk workers were asked to select the best paraphrase given an input text and the output of two models, based on three criteria: fluency, meaning, and dissimilarity. Figure 1 shows the user interface used for the human evaluation in the original study. The verbatim instructions of the task included in the user interface are provided below:

- *Which system output is the most **fluent and grammatical**?*

- *To what extent is the **meaning** expressed in the original sentence **preserved** in the rewritten version, with no additional information added?*

- *Does the rewritten version use **different words or phrasing** to the original? You should choose the system that uses the most different words or word order.*

The authors provided additional information regarding the human evaluation in the appendix of their paper. Importantly, they reported utilizing MTurk's feature to make HITs available only in specific regions, setting their region availability to the United States and the United Kingdom. Furthermore, they reported that participants were compensated for their time at a rate above the living wage in the regions selected.

Ultimately, in comparing paraphrase pairs the authors evaluated 300 sentences sampled equally from the three datasets, with paraphrases generated by each model resulting in a total of 1800 paraphrases.[5] For a particular pair of two system outputs for a given input sentence, separately for each of the three criteria, a given system received +1 or -1 depending on whether it was chosen as the best (+1) or worst (-1). The final scores for each model were then calculated by averaging the scores across all of that model's scored samples for a particular criterion. This scoring process is

known as Best-Worst Scaling (Louviere and Woodworth, 1991; Louviere et al., 2015).

According to the authors, HRQ-VAE was found to be *more fluent* and *more diverse* while maintaining a *similar meaning* to the original sentence. Figure 4 in their paper shows the results of the human evaluation. We identified five unique **claims** based on the human evaluation results in the original paper:

- **Claim 1:** The VAE baseline is the best at preserving meaning.

- **Claim 2:** The VAE baseline is the worst at introducing variation to the output.

- **Claim 3:** HRQ-VAE better preserves the original intent compared to the other systems.

- **Claim 4:** HRQ-VAE introduces more diversity than VAE.

- **Claim 5:** HRQ-VAE generates much more fluent output than VAE.

## 3.2. Scope of Reproduction

Our goal was to repeat the allocated experiment as closely as possible to the original study. We focused on a narrow scope of the original paper: we sought to reproduce the outcomes of the human evaluation experiments for the *meaning* criterion. We set up the experiment using all information available to us from the original paper (Hosking et al., 2022) and from follow-up communications with the authors by the ReproHum leadership team. We filled Human Evaluation Datasheet (HEDS) (Shimorina and Belz, 2022) containing the details of the human evaluation experiment. The HEDS is released in the ReproNLP central GitHub repository for HEDS documents.[6]

## 3.3. Additional Information Obtained from the Original Authors

While we did not directly communicate with the original authors, the ReproHum team provided us with additional information obtained from them. Specifically, the authors shared the exact outputs that they evaluated and the user interface that they used for the human evaluation. Crucially, the authors noted that they used *attention checks* (control samples with known labels). Each task contained two control samples; in one control sample, the system was a "distractor" and the output was a random sample with a completely different meaning that should clearly never be chosen as best for the *meaning* criterion. The other control sample was

---

[4]Note that *HIT* is a term used on MTurk to refer to a single task or job that a worker can complete; we use the terms *task* and HIT interchangeably in this paper.

[5]There were four systems; for each comparison, we selected two out of four: $\binom{4}{2} = 6$. With the resulting six unique comparisons for each of the 300 sentences, we have a total of $6 \times 300 = 1800$ comparisons.

[6]https://github.com/nlp-heds/repronlp2024

when the system's output was the same as the input, which should clearly never be chosen as best for the *dissimilarity* criterion. Note that the second control sample was not relevant to our reproduction, as we were reproducing the results for the *meaning* criterion. In their communication, the authors mentioned that HITs for which either of these attention checks were failed were rejected and resubmitted to MTurk. Additionally, they reported compensating participants with $3.50 per HIT with an expected completion time of 20 minutes.

## 3.4. Notes on Experimental Design

The original design of the human evaluation did not consider cases in which both outputs were equally good (a tie). Although we would have preferred to include this option, we followed the original design. Moreover, in analyzing the outputs we uncovered a slight imbalance in the number of samples selected from each dataset. Specifically, while QQP had 100 samples, MSCOCO had 102 samples and Paralex had 98 samples.

## 3.5. Known Deviations from the Original Experiment

We are aware of several deviations in our reproduction from the original experiment, and we detail these below. We do not believe that these deviations had a major impact on our reproduction results.

**Crowdsourcing Platform:** Our biggest deviation from the original experiment was in the crowdsourcing platform used. While the original study had utilized MTurk, we used Prolific.[7] This decision was made across all experiments in the ReproHum project to ensure consistency, due to limitations in credit usage on MTurk and the administrative overhead of managing the funds for different experiments.

Prolific survey design is different from MTurk, and we had to adapt the original survey design to the Prolific platform. To be more specific, setting up a survey similar to the structure of HITs was only possible using external survey tools. The ReproHum team shared the code for hosting a server to run the survey. We used a modified version of the code with additional checks to ensure the validity of the responses. Furthermore, we added thread safety to prevent race conditions, where two or more threads try to access or modify the same data at the same time, leading to unpredictable or incorrect results.

**Region Control:** Our reproduction also deviated slightly in terms of participant region control. While the original authors had limited their HIT availability regions to the United States and the United Kingdom, we followed the region control guidelines of all experiments in the ReproHum project. This meant that participants from Australia and Canada were also included in addition to the United States and the United Kingdom.

**Participant Selection:** The authors reported filtering participants with approval rates less than 96%, and required that participants had completed at least 5000 HITs. In contrast, we set the approval rate to 99% and the minimum number of HITs completed to 200. This decision was based on the recommendations from Prolific to ensure high-quality participants.[8]

**Failed Attention Checks:** The original authors reported rejecting HITs for which the attention checks were failed. We did not reject any HITs based on attention checks per recommendations from the ReproHum team; however, we solicited new responses for tasks that failed attention checks.

**Participation Limit:** The original paper did not report whether a participant could respond to multiple HITs; we assume that no controls were in place for this. In Prolific, participants cannot respond to the same study more than once, even though the input data may be different.

**Expected Completion Time:** The original authors reported that the expected completion time for a HIT was 20 minutes. Our survey differed from the original study since we only collected responses for the *meaning* criterion. We ran several surveys to estimate the time it would take to complete the task. Ultimately, we set the expected completion time to 8 minutes.

**Payment:** The original authors reported compensating participants with $3.50 per HIT for 20 minutes of work, resulting in an hourly rate of $10.50. We followed the guidelines of the ReproHum project, setting the wage as the minimum living wage in the United Kingdom (which was higher than our local minimum wage). At the time of data collection, this value was £12 which was equivalent to $15.14 using the exchange rate between UK and US currency at that time. To be more specific,

---

[7]https://www.prolific.com/

[8]https://www.prolific.com/resources/find-fil ter-favourite-how-to-select-participants-for-a i-tasks

the participants received £1.60 or $2 for 8 minutes of work.

**User Interface:** Our institutional consent forms were required to be much more detailed than those used in the original study, and this was beyond our control. To ensure that the participants were not overwhelmed, we split the welcome, instructions, and task into three separate pages. We have included images of the user interface used for the reproduction in the appendix (Figures 3, 4, and 5).

**Data Analysis:** The source code for the data analysis was intentionally left out and we were asked to write our own code to analyze the data. We also conducted additional analyses to better understand the data. We report our findings from these analyses in §5.

## 4. Quantified Reproducibility Assessment

We followed the standardized procedure for reproducibility assessment as outlined by the ReproHum team. For single numerical result scores, we calculated the coefficient of variation (CV) to quantify the precision of the results. The CV is calculated as the ratio of the standard deviation of the results to their mean. It serves as a measure of relative variability, and it is useful for comparing the precision of different experiments. We adjusted the CV for small sample sizes as reported by Belz (2022), and refer to this adjusted CV using the notation *CV\**. Furthermore, the results are shifted by 100 to ensure the mean is positive, as the original scores were in the range of -100 to 100.

For sets of numerical scores, we calculated Pearson and Spearman correlations between the reproduced and original results. The Pearson correlation measures the linear relationship between two sets of scores, and the Spearman correlation measures the monotonic relationship between two sets of scores. Both correlations range from -1 to 1, where 1 indicates a perfect positive correlation, -1 indicates a perfect negative correlation (suggesting that the outcomes are diametrically opposed), and 0 indicates no correlation. Using these metrics, we assessed how closely the reproduced results aligned with the original results.

## 5. Results

### 5.1. Study Analysis

According to the summary statistics provided by Prolific, the median time spent on the survey was 7 minutes and 12 seconds. With this time, the actual hourly rate was calculated to be £13.30. With filters

| System | Orig | Ours | CV* | $r$ | $\rho$ |
|---|---|---|---|---|---|
| VAE | 36 | 37.04 | 0.76 | | |
| Latent BoW | -16 | -14.52 | 1.74 | | |
| Separator | -24 | -29.78 | 7.88 | 0.99 | 1 |
| HRQ-VAE | 4 | 7.26 | 3.08 | | |

Table 1: Overview comparing the original and reproduced versions of the human evaluation, including precision metrics to reflect the degree of reproducibility. Pearson's correlation is represented by $r$ and Spearman's correlation is represented by $\rho$. CV* is computed using $n$=2. *Orig* refers to the original results reported by Hosking et al. (2022).

| Sys. | Win # | Loss # | Best-Worst Score | Best-Worst Scale | Win % |
|---|---|---|---|---|---|
| VAE | 1850 | 850 | 1000 | 37.04 | 68.52 |
| Lat. BoW | 1154 | 1546 | -392 | -14.52 | 42.74 |
| Sep. | 948 | 1752 | -804 | -29.78 | 35.11 |
| HRQ-VAE | 1448 | 1252 | 196 | 7.26 | 53.63 |

Table 2: Additional details from our own reproduced human evaluation. *Lat. BoW* refers to the *Latent Bag-of-Words* system, and *Sep.* refers to the *Separator* system.

set for region control and acceptance rate, 51,430 of 152,649 possible participants were eligible to participate in the study; our 180 participants were selected from this pool. We had to repeat one task due to failed attention checks, making the total number of participants $n$=181.

Aside from data available in Prolific, we collected additional data from the survey. Particularly, we collected the time spent on each page of the survey. We present the histogram of time spent on each page of the survey in Figure 6 (in Appendix B). Furthermore, we present the empirical cumulative distribution function (eCDF) of the time spent on each page of the survey in Figure 7 (in Appendix B). Note that this data may not be entirely reliable as participants were given an hour to complete the survey, and the time spent on each page was not necessarily indicative of the time spent on the task (e.g., participants may have stepped away from the computer while leaving the page open). Nonetheless, we consider it a reasonable proxy for the time spent on the task.

The 50th percentiles (median) of the time spent on the welcome, instructions, and task pages were 13, 53, and 328 seconds, respectively. Additionally, we observed that the 90th percentiles of the time spent on the welcome and instruction pages were

82 and 92 seconds, respectively. In other words, the eCDF suggests that 90% of the participants spent less than 82 seconds on the welcome page and 92 seconds on the instruction page. The task page eCDF suggests that the 80% percentile of time spent on the task page was 434 seconds, meaning that 80% of the participants spent less than 434 seconds on the task page. Recall that the total time allotted for the survey was 480 seconds (8 minutes).

## 5.2. Reproduction Results

Table 1 shows the results of the human evaluation for the selected criterion, comparing the outcomes from the original and reproduced experiments. Overall, we observe that our results are very close to the scores originally reported (Hosking et al., 2022). This is reflected in low CV* values for all the systems. Pearson correlation and p-value are $r$=0.99 and $p$=0.01, respectively. Similarly, Spearman correlation and p-value are $\rho$=1.00 and $p$=0.00. Both Pearson and Spearman correlations are very high, indicating a strong relationship between the original and reproduced scores. Figure 2 presents this same information in the format used by the original paper, showing best-worst scaling outcomes for the four systems compared in the original paper and in our reproduction.

In Table 2, we include additional details from our own reproduced human evaluation. We report the number of wins and losses for each system, the best-worst score outcome (the sum of all scores of +1 or -1 that the system received), and the best-worst scale outcome. We also report the percentage of wins for each system. We used Krippendorff's alpha to evaluate the agreement among the categorical responses collected, resulting in a value of $\alpha$=0.51. This metric was not included in the original study, preventing a direct comparison of our findings.

For statistical analysis, we employed ANOVA to determine significant differences among the means of multiple independent groups. We measured effect size using partial eta squared ($\eta^2$), which yielded a large effect size of 0.17 for the ANOVA test. With a sample size of 300 and $\alpha$=0.05, the calculated test power was 0.67, falling below the recommended threshold of 0.80. Achieving a power of 0.80 would require a sample size of 395. In conducting the ANOVA test, we observed an F value of 79.93 with a corresponding $p$=3.97e-47. Subsequently, we used Tukey's HSD test to identify significant differences between individual groups, revealing significant distinctions among all groups.

Overall, given our reproduced results' similarity to and correlation with the originally reported results, we could easily confirm two out of five of the original claims based on the human evaluation re-

| Claim | Verification |
|---|---|
| The VAE baseline is the best at preserving meaning. | Verified |
| The VAE baseline is the worst at introducing variation to the output. | Out of Scope |
| HRQ-VAE better preserves the original intent compared to the other systems. | Verified |
| HRQ-VAE introduces more diversity than VAE. | Out of Scope |
| HRQ-VAE generates much more fluent output than VAE. | Out of Scope |

Table 3: Claims and verifications.

sults. The other three claims were out of scope for our reproduction, as they pertained to criteria other than *meaning*. We summarize the claims and our verification in Table 3.

## 6. Discussion

Since we found this experiment to be underpowered, combining the data collected in our reproduction with the parallel work of the ReproHum project could provide a more robust analysis. This would allow us to draw more reliable conclusions about the reproducibility of the original study. Nonetheless, if the results of the other reproduction are consistent with ours, we believe this experiment is a good candidate for the next round of the ReproHum project, where some variations could be introduced to further investigate the reproducibility of the original study. Replacing the attention check that is not relevant to the *meaning* criterion with a more relevant one could be a good starting point.

We followed the Prolific recommended guidelines for selecting participants, setting the approval rate to 99% and the minimum number of accepted tasks to 200. The problem with this approach is that the number of accepted tasks inflates over time. A better alternate approach would be to select the top $k$% of workers based on the total number of accepted tasks. A similar concern was raised by González Corbelle et al. (2023). Considering that data collection is essential to machine learning and NLP research, it is important to ensure the quality of the data collected. Lastly, we observed that some work was submitted in timezones other than those associated with the regions selected. This could be due to participants using VPNs or other methods to change their location. In general, timezones are not reliable and can be easily changed. Thus, this is a complex issue that re-
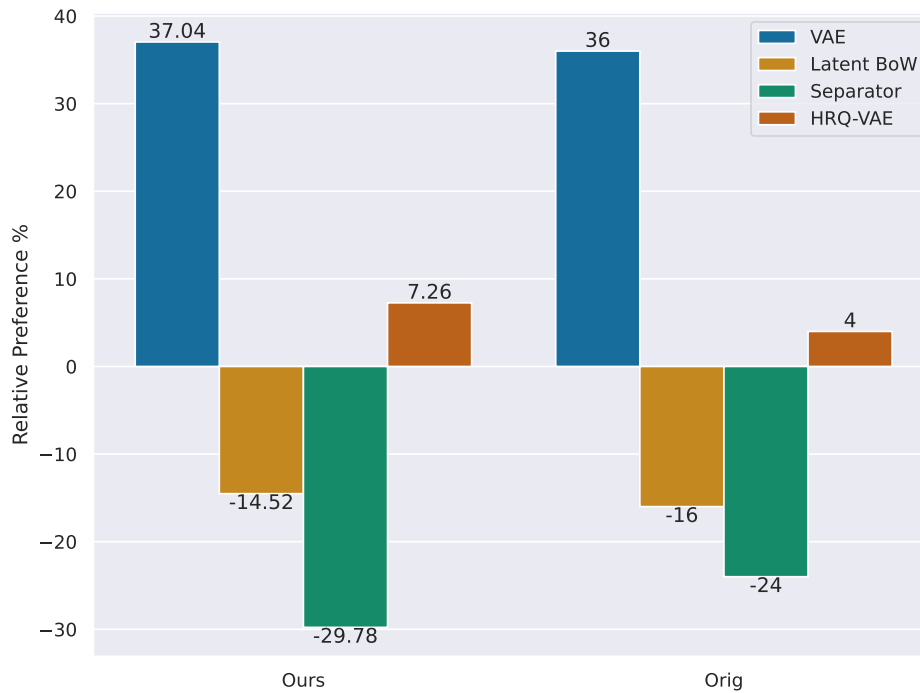
Figure 2: Results of the human evaluation, comparing the original and reproduced systems. Results are presented in the same format used in the original paper.

quires action, cooperation, and transparency from crowdsourcing platforms to ensure the quality of the data collected.

Finally, Platek et al. (2023) report having difficulties setting up the user interface for their reproduction. They suggest utilizing a Docker image containing all the dependencies. We believe that this is a good practice. Considering that our server setup for ReproHum reproductions is customized and unique, we have included the docker compose configuration to bring up the server with all the dependencies and tasks in a separate repository.[9]

## 7. Conclusion

In this reproduction, we studied the extent to which the human evaluation reported in "Hierarchical Sketch Induction for Paraphrase Generation" is reproducible, narrowing our scope to a single evaluation criterion (*meaning*). We systematically and carefully reproduced the experiment as reported in the original paper to ensure consistency with the original settings to the extent possible. Through a comparison of our reproduced results with those achieved in the original paper using CV*, Pearson's correlation, and Spearman's correlation, we believe that the human evaluation conducted by the original authors has a high degree of reproducibil-

ity. This reflects the quality of the design of the experiment. This work would not have been possible without the support of the ReproHum project and the original authors. We hope that our work will contribute to ongoing efforts to improve the reproducibility of research in the field of NLP.

## Acknowledgments

## References

Mohammad Arvan and Natalie Parde. 2023. Human evaluation reproduction report for data-to-text generation with macro planning. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 89–96, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Anya Belz. 2022. A metrological perspective on reproducibility in NLP. *Comput. Linguistics*, 48(4):1125–1135.

---

[9] https://github.com/mo-arvan/reprohum-prolific-webapp

Anya Belz and Craig Thomson. 2023. The 2023 ReproNLP shared task on reproducibility of evaluations in NLP: Overview and results. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 35–48, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Jackie Cheung, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Huiyuan Lai, Chris van der Lee, Emiel van Miltenburg, Yiru Li, Saad Mahamood, Margot Mieskes, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Pablo Mosteiro Romero, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *The Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew M. Dai, Rafal Józefowicz, and Samy Bengio. 2016. Generating sentences from a continuous space. In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, CoNLL 2016, Berlin, Germany, August 11-12, 2016*, pages 10–21. ACL.

Timothy M Errington, Alexandria Denis, Nicole Perfito, Elizabeth Iorns, and Brian A Nosek. 2021a. Challenges for assessing replicability in preclinical cancer biology. *elife*, 10:e67995.

Timothy M Errington, Maya Mathur, Courtney K Soderberg, Alexandria Denis, Nicole Perfito, Elizabeth Iorns, and Brian A Nosek. 2021b. Investigating the replicability of preclinical cancer biology. *Elife*, 10:e71601.

Anthony Fader, Luke Zettlemoyer, and Oren Etzioni. 2013. Paraphrase-driven learning for open question answering. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, ACL 2013, 4-9 August 2013, Sofia, Bulgaria, Volume 1: Long Papers*, pages 1608–1618. The Association for Computer Linguistics.

Yao Fu, Yansong Feng, and John P. Cunningham. 2019. Paraphrase generation with latent bag of words. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 13623–13634.

Mingqi Gao, Jie Ruan, and Xiaojun Wan. 2023. A reproduction study of the human evaluation of role-oriented dialogue summarization models. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 124–129, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Javier González Corbelle, Jose Alonso, and Alberto Bugarín-Diz. 2023. Some lessons learned reproducing human evaluation of a data-to-text system. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 49–68, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Tom Hosking and Mirella Lapata. 2021. Factorising meaning and form for intent-preserving paraphrasing. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 1405–1418. Association for Computational Linguistics.

Tom Hosking, Hao Tang, and Mirella Lapata. 2022. Hierarchical sketch induction for paraphrase generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 2489–2501. Association for Computational Linguistics.

David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In *Proceedings of the 13th International Conference on Natural Language Generation, INLG 2020, Dublin, Ireland, December 15-18, 2020*, pages 169–182. Association for Computational Linguistics.

Manuela Hürlimann and Mark Cieliebak. 2023. Reproducing a comparative evaluation of German text-to-speech systems. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 136–144, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Takumi Ito, Qixiang Fang, Pablo Mosteiro, Albert Gatt, and Kees van Deemter. 2023. Challenges in reproducing human evaluation results for role-oriented dialogue summarization. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 97–123, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Filip Klubička and John D. Kelleher. 2023. HumEval'23 reproduction report for paper 0040: Human evaluation of automatically detected over- and undertranslations. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 153–189, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Yiru Li, Huiyuan Lai, Antonio Toral, and Malvina Nissim. 2023. Same trends, different answers: Insights from a replication study of human plausibility judgments on narrative continuations. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 190–203, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: common objects in context. In *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V*, volume 8693 of *Lecture Notes in Computer Science*, pages 740–755. Springer.

Jordan J Louviere, Terry N Flynn, and Anthony Alfred John Marley. 2015. *Best-worst scaling: Theory, methods and applications*. Cambridge University Press.

Jordan J Louviere and George G Woodworth. 1991. Best-worst scaling: A model for the largest difference judgments. Technical report, Working paper.

Saad Mahamood. 2023. Reproduction of human evaluations in: "it's not rocket science: Interpreting figurative language in narratives". In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 204–209, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Margot Mieskes and Jacob Georg Benz. 2023. h_da@ReproHumn – reproduction of human evaluation and technical pipeline. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 130–135, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Open Science Collaboration. 2015. Estimating the reproducibility of psychological science. *Science*, 349(6251):aac4716.

Ondrej Platek, Mateusz Lango, and Ondrej Dusek. 2023. With a little help from the authors: Reproducing human evaluation of an MT error detector. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 145–152, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Hong Sun and Ming Zhou. 2012. Joint learning of a dual SMT system for paraphrase generation. In *The 50th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference, July 8-14, 2012, Jeju Island, Korea - Volume 2: Short Papers*, pages 38–42. The Association for Computer Linguistics.

Emiel van Miltenburg, Anouck Braggaar, Nadine Braun, Debby Damen, Martijn Goudbeek, Chris van der Lee, Frédéric Tomas, and Emiel Krahmer. 2023. How reproducible is best-worst scaling for human evaluation? a reproduction of 'data-to-text generation with macro planning'. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 75–88, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Lewis Watson and Dimitra Gkatzia. 2023. Unveiling NLG human-evaluation reproducibility: Lessons learned and key insights from participating in the ReproNLP challenge. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 69–74, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

## A. Reproduction User Interface

We show our reproduced interface for the human subject consent page for the human evaluation in Figure 3. Participants were required to consent by clicking the "Accept & Continue" link prior to taking part in the evaluation. In Figures 4 and 5 we present the participant and template views for the reproduced evaluation, respectively.

## B. Time Spent on Survey

In Figures 6 and 7 we report the amount of time spent by participants on the reproduced evaluation. Time was recorded for each page of the survey. Figure 6 shows a histogram of the number
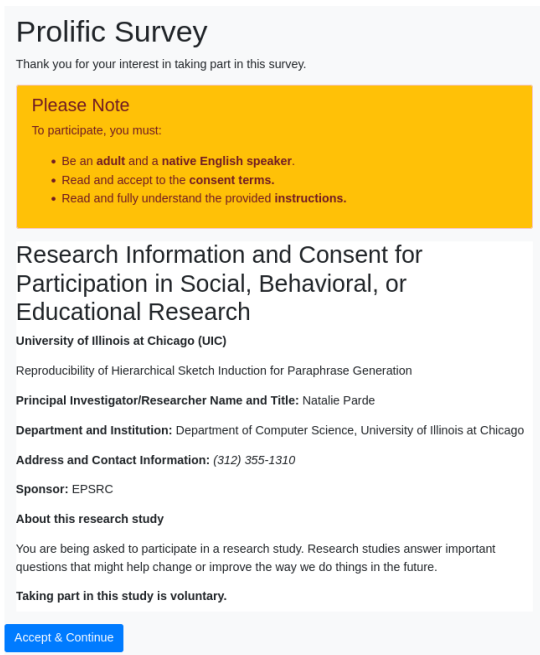
Figure 3: Reproduced interface for the human evaluation (consent page).

of seconds spent on each page, whereas Figure 7 computes and displays an empirical cumulative distribution function for this data.
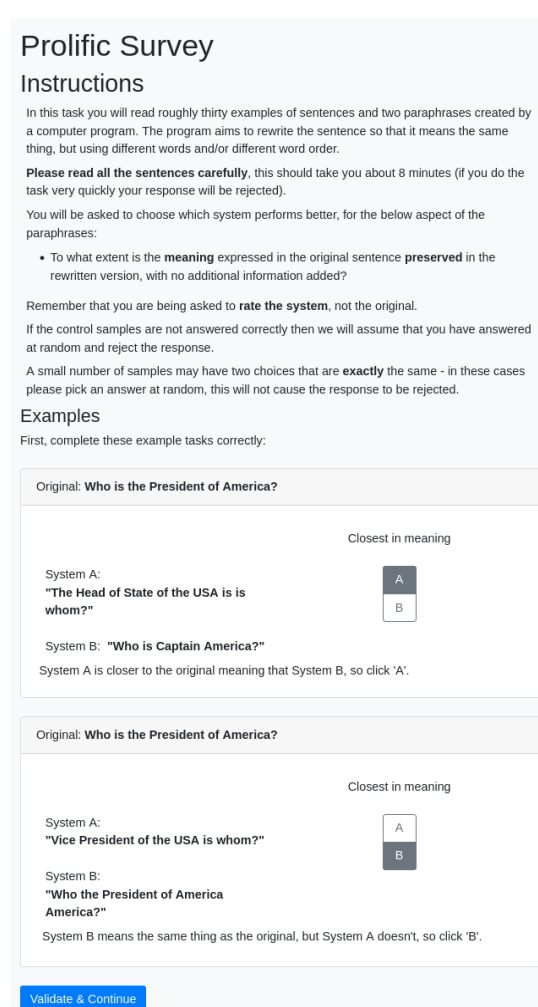


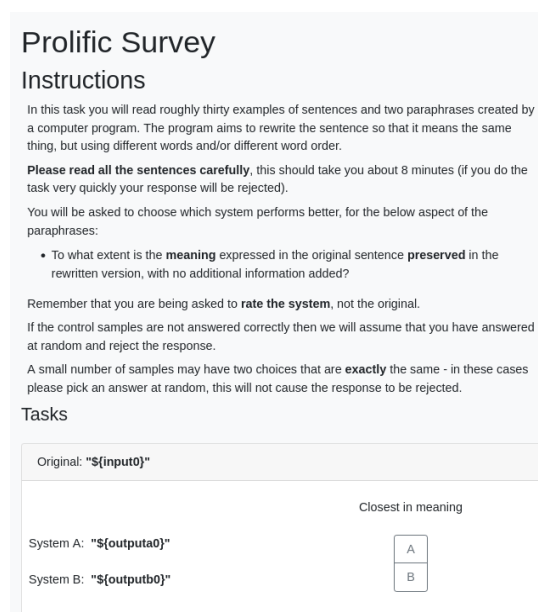Figure 4: Reproduced interface for the human evaluation (participant view).



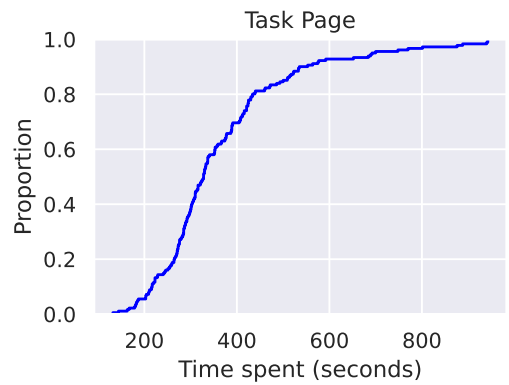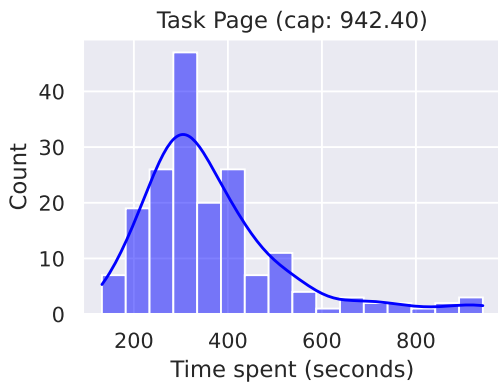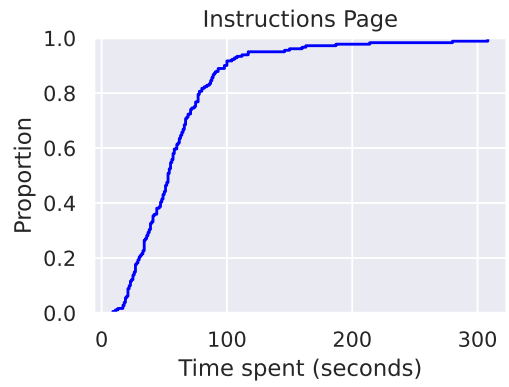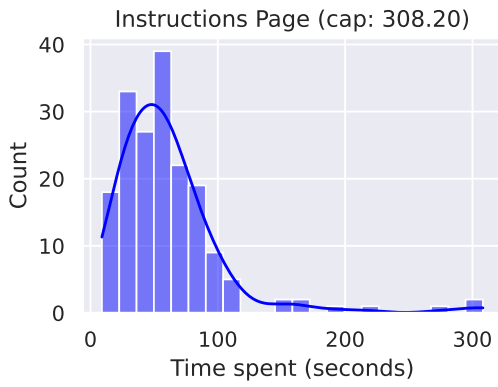Figure 5: Reproduced interface for the human evaluation (template view).

Figure 6: Histogram of seconds spent on each page of the survey. Note that each histogram is capped to ensure readability.



Figure 7: Empirical Cumulative Distribution Function (eCDF) of seconds spent on each page of the survey.

# ReproHum #0712-01: Reproducing Human Evaluation of Meaning Preservation in Paraphrase Generation

**Lewis Watson, Dimitra Gkatzia**

Edinburgh Napier University

{L.Watson, D.Gkatzia}@napier.ac.uk

## Abstract

Reproducibility is a cornerstone of scientific research, ensuring the reliability and generalisability of findings. The ReproNLP Shared Task on Reproducibility of Evaluations in NLP aims to assess the reproducibility of human evaluation studies. This paper presents a reproduction study of the human evaluation experiment in "Hierarchical Sketch Induction for Paraphrase Generation" by Hosking et al. (2022). The original study employed a human evaluation on Amazon Mechanical Turk, assessing the quality of paraphrases generated by their proposed model using three criteria: meaning preservation, fluency, and dissimilarity. In our reproduction study, we focus on the meaning preservation criterion and utilise the Prolific platform for participant recruitment, following the ReproNLP challenge's common approach to reproduction. We discuss the methodology, results, and implications of our reproduction study, comparing them to the original findings. Our findings contribute to the understanding of reproducibility in NLP research and highlights the potential impact of platform changes and evaluation criteria on the reproducibility of human evaluation studies.

**Keywords:** reproducibility, NLG, paraphrase generation, human evaluation

## 1. Introduction

Reproducibility is a fundamental principle of scientific research, ensuring that findings can be independently verified and built upon by the wider research community. In the field of Natural Language Generation (NLG), reproducibility is particularly challenging due to the complex nature of the tasks (Belz et al., 2023) and the use of human assessments for the evaluation of NLG approaches (Gehrmann et al., 2023; Howcroft et al., 2020). Recently, the reproducibility of NLP studies has been called into question, with concerns raised about the reliability and generalisability of reported findings (Belz et al., 2021).

**The ReproNLP Challenge**   To address the issue of reproducibility in NLP, the ReproNLP/ReproGen challenge was established to assess the reproducibility of human evaluation studies (Belz et al., 2020), under three conditions: (1) reproduction of evaluation results of pre-selected papers based on information of the original paper and additional information by the authors; (2) reproduction of evaluation results by the same authors, i.e. own study; (3) reproduction of a pre-selected study using information provided by the ReproNLP organisers only (Belz and Thomson, 2023).

The 2023 round of reproduction studies provided a wealth of lessons learnt. The evaluators' background and qualifications were identified as important factors in obtaining consistent results as discrepancies in these can lead to varying results. (González Corbelle et al., 2023; Watson and Gkatzia, 2023; Mieskes and Benz, 2023; Li et al.,

2023; Mahamood, 2023). The number of ratings obtained per item and worker are also important for obtaining statistically similar results (van Miltenburg et al., 2023; Ito et al., 2023; Gao et al., 2023). In addition, Ito et al. (2023) highlight that errors in statistical analyses can prohibit reproducibility. Technical issues can prohibit replications of studies that can be overcome through the use of Docker (Platek et al., 2023) and the provision of code used for crowdsourcing and the analysis of results (Mahamood, 2023). However, Klubička and Kelleher (2023) used a different user interface for their reproduction study than the original authors and were able to confirm the results of the original study. Discrepancies in the study design have also been identified as an issue in reproducibility (Platek et al., 2023; Gao et al., 2023), while Platek et al. (2023) advocate that setups with a minimal range of potential answers, particularly those with binary questions, are simpler to duplicate and should be favoured over more intricate setups whenever feasible.

In the 2023 round, we reproduced the human evaluation as close as possible to the methodology used by the original authors (Watson and Gkatzia, 2023). In the 2024 round, we experimented with a platform change, Prolific instead of Amazon Mechanical Turk, and we focused on only one quality criterion (meaning preservation) as outlined in the 2024 challenge design (Belz and Thomson, 2024).

**Hierarchical Sketch Induction for Paraphrase Generation**   In this paper, we focus on reproducing a single quality criterion from the human evaluation component from the study "Hierarchical Sketch

Induction for Paraphrase Generation" by Hosking et al. (2022). The original study proposed a novel approach to paraphrase generation using hierarchical sketch induction and conducted a human evaluation on Amazon Mechanical Turk to assess the quality of the generated paraphrases based on three criteria: meaning preservation, fluency, and dissimilarity.

**Human Evaluation Datasheet (HEDS)**  As part of our reproduction study, we have completed the Human Evaluation Datasheet (HEDS) (Shimorina and Belz, 2022), a standardised template for documenting human evaluation experiments in NLP. The HEDS framework aims to promote reproducibility and facilitate meta-evaluation of evaluation methods by providing a consistent format for recording the details of human evaluations. Our completed HEDS document is available on our project's GitHub repository[1] and has also been contributed to the central HEDS repository maintained by the ReproNLP organisers[2]. This central repository serves as a comprehensive resource for HEDS documents from all participating teams, enabling access and comparison of human evaluation methodologies across different studies. By adhering to the HEDS framework and sharing our documentation, we aim to support the broader goal of improving the reliability and generalisability of human evaluation practices in the field.

## 1.1. Objectives and Hypotheses

The main objective of our reproduction study is to assess the reproducibility of the human evaluation results reported in Hosking et al. (2022). We aim to answer the following research questions:

1. To what extent can the human evaluation results be reproduced using a different participant recruitment platform (Prolific instead of Amazon Mechanical Turk)?

2. How does focusing on a single evaluation criterion (meaning preservation) affect the reproducibility of the results compared to the original study, which used three criteria?

Based on these research questions, we hypothesise that:

1. The change in participant recruitment platform may lead to some differences in the evaluation results, but the overall trends should remain consistent with the original study.

2. Focusing on a single evaluation criterion may result in higher reproducibility compared to the

original study, as it reduces the complexity of the task and the potential for variability in participant judgements.

## 2. Original Study

### 2.1. Methodology

The original study by Hosking et al. (2022) proposed a novel approach to paraphrase generation called Hierarchical Refinement Quantized Variational Autoencoders (HRQ-VAE). The HRQ-VAE model learns to generate paraphrases by first inducing a syntactic sketch of the input sentence, which captures its syntactic structure at varying levels of granularity. The model then generates the final paraphrase based on the induced sketch and the original sentence's meaning representation.

To evaluate the quality of the generated paraphrases, the authors conducted a human evaluation study on Amazon Mechanical Turk (AMT)[3]. The annotators were required to have an approval rate of >96%, be located in the United States or United Kingdom, and have completed >5000 HITs, workers were paid $3.50USD/hr. They compared the HRQ-VAE model's output to paraphrases generated by three other baseline models, namely, Gaussian Variational AutoEncoder (VAE Bowman et al. 2016), Separator (Hosking and Lapata, 2021) and Latent bag-of-words (BoW, Fu et al. 2019)[4].

The human evaluation tasks were created using 300 input sentences sampled equally from three datasets: Paralex (Fader et al., 2013), Quora Question Pairs (QQP) (Chen et al., 2017), and MSCOCO (Lin et al., 2014). For each input sentence, the paraphrases generated by the HRQ-VAE model and the baseline models were presented to the AMT workers, who were asked to rate the paraphrases based on three criteria:

1. **Meaning preservation**: The extent to which the generated paraphrase preserves the meaning of the original input sentence.

2. **Fluency**: The fluency and grammaticality of the generated paraphrase.

3. **Dissimilarity**: The degree to which the generated paraphrase differs from the original input sentence in terms of word choice and sentence structure.

Each comparison was evaluated by 3 distinct AMT workers, resulting in a total of 900 judgements (300 sentences × 3 judgements per sen-

---

tence). Each task contained 32 paraphrase questions, including 2 attention checks.

The first attention check focused on the meaning criteria and consisted of comparisons where one paraphrase is generated by a "distractor" model designed to produce output with a completely different meaning. The second attention check focused on the dissimilarity criteria where the paraphrase would be the same as the input. Where a participant failed the attention check, their results were discarded.

## 2.2. Results

The original study reported the human evaluation results as relative preference scores for each of the three dimensions (meaning, dissimilarity, and fluency) across the four models: HRQ-VAE, Separator, Latent BoW, and VAE. The relative preference scores were calculated by assigning a score of +1 when a system was selected, -1 when the other system was selected, and taking the mean over all samples.

Key findings from the original study include:

- The VAE baseline achieved the highest relative preference score for meaning preservation (+36%) but the lowest for dissimilarity (-33%), indicating that while it best preserved the original sentence's meaning, it introduced the least variation in the generated paraphrases.

- The HRQ-VAE model offered the best balance between meaning preservation (+4%) and dissimilarity (-3%), demonstrating its ability to generate paraphrases that maintain the original meaning while introducing diversity.

- In terms of fluency, the HRQ-VAE model outperformed Separator and Latent BoW, with a relative preference score of +8%.

These findings highlighted the effectiveness of the proposed hierarchical sketch induction approach in generating high-quality paraphrases that strike a balance between meaning preservation and dissimilarity while maintaining fluency.

# 3. Reproduction

## 3.1. Methodology

Our reproduction study aims to assess the reproducibility of the human evaluation results reported in the original study by Hosking et al. (2022). We follow the ReproNLP challenge's common approach to reproduction (Belz et al., 2020), with some modifications to the participant recruitment process and the evaluation criteria.

### 3.1.1. Participant Recruitment

We recruited participants using the Prolific crowdsourcing platform[5], which differs from the original study's use of Amazon Mechanical Turk (AMT). Participants were sourced from the United Kingdom, Canada, the United States, and Australia to ensure a diverse sample and adhere to the ReproNLP Challenge. To prevent overlap with the participant pool of another lab conducting a similar reproduction study, we exclude participants who have taken part in their study. Additionally, in accordance with the ReproHum common procedure for calculating fair pay (Belz et al., 2023), participants were paid £2. This was calculated by assuming the reduced complexity task should take around 10 minutes and paying £12/hr. The median time to complete the task was 8 minutes and our average reward per hour came to £14.75.

In contrast to the original study, we did not impose any restrictions on the participants' approval rate or number of previously completed tasks on Prolific.

### 3.1.2. Evaluation Tasks and Procedure

We use the same set of 300 sentences as in the original study. These sentences are divided into 60 distinct tasks (each needing three participant ratings, therefore requiring 180 participants), each containing 32 paraphrase questions, including 2 attention checks.

A single question in a task consisted of an original sentence along with two corresponding paraphrases, each generated by distinct models. Contrasting with the methodology of the original study, our reproduction concentrated solely on a singular criterion. This decision was informed by the preliminary ReproHum findings, which indicated that tasks of lower complexity yielded enhanced reproducibility (Belz et al., 2023). The participants' assigned task was to identify the paraphrase that most effectively retained the meaning of the original sentence.

Each distinct task was evaluated by 3 participants, resulting in a total of 180 participant results (60 distinct tasks × 3 participants per task). After removing the attention check questions, we obtain a total of 1,800 final average comparisons (5760 total evaluations ÷ 3 participants = 1920 average from participants, then 1920 - 120 attention checks = 1800 final). The four models being evaluated in this study are the same as in the original study: VAE, Latent BoW, Separator, and HRQ-VAE.

## 3.2. Attention Check

To ensure the quality of the collected data, we incorporate an attention check mechanism in our

223

---

reproduction study, following the same approach as the original paper. The attention check consists of comparisons where one of the paraphrases is generated by a "distractor" model, which is designed to produce output with a completely different meaning from the original sentence. If a participant selected the distractor model, their responses were discarded and reran but we did still pay the participant. We had 5 failed attention checks in the initial run of the reproduction, and then a further 1 failed attention checks on the rerun totalling 6 failed attention checks. There are 2 attention checks per task, and with 60 distinct tasks, there are a total of 120 attention check questions (2 x 60 = 120).

We decided to include the original study's second attention check question, to minimise the differences between the original study and the reproduction however, the data was not used for analysis.

### 3.3. Preference Calculation

To analyse the results of the reproduction study, we follow the same approach as the original study. For each comparison between two paraphrases, we assign a score of +1 to the model whose paraphrase is selected by the participant as better preserving the meaning of the original sentence. Conversely, the model whose paraphrase is not selected receives a score of -1. This scoring method allows us to calculate the relative preference for each model. The analysis is performed using a Python script, which can be found alongside our raw results on GitHub[6]. The script reads the data from a CSV file and iterates over each unique task number (1-60). For each task, it examines the participant responses for the meaning preservation criterion across all 32 comparisons, excluding the attention check questions.

For each comparison, the script determines the preferred model based on the majority vote across the three participants. If model A is preferred, it receives a score of +1, while model B receives a score of -1, and vice versa. These scores are accumulated for each model across all comparisons.

After processing all the comparisons, the script calculates the total number of comparisons (1800 once the attention checks have been removed) and the average number of preferences across all models. Finally, it computes the relative preference percentage for each model by dividing its accumulated score by the total number of comparisons and multiplying by 100.

The resulting relative preference percentages provide insights into the performance of each model in terms of meaning preservation, as judged by the participants in the reproduction study. These percentages can then be compared to the original

study's results to assess the reproducibility of the findings.

### 3.4. Differences from the Original Study

Our reproduction study differs from the original study in the following aspects:

- We use the Prolific platform for participant recruitment instead of Amazon Mechanical Turk.

- We do not impose restrictions on participants' approval rate or number of previously completed tasks.

- We focus on a single evaluation criterion (meaning preservation) instead of three criteria (meaning, dissimilarity, and fluency).

- We recruited participants from the United Kingdom, Canada, the United States, and Australia as opposed to just the UK and USA like the original study.

These differences allow us to investigate the impact of participant recruitment platforms, screening criteria, and evaluation criteria on the reproducibility of the human evaluation results. It is important to note that some of these changes were planned, such as focusing on a single evaluation criterion and recruiting participants from additional countries, while others, like the omission of participant approval rate and task completion restrictions, were unintentional.

The omission of Prolific filters was an oversight, however it highlights the challenges of conducting reproduction studies with complete accuracy. As Thomson et al. (2024) argues, mistakes might occur in many human evaluations, and there is no evidence to suggest that all published studies are entirely mistake-free. Despite our best efforts to adhere to the original study's methodology, this unintended difference in participant screening criteria may have introduced additional variability in our reproduction results.

## 4. Results

In this section, we present the results of our reproduction study and compare them with the findings of the original study by Hosking et al. (2022).

Figure 1 illustrates the relative preference results from our reproduction study. The HRQ-VAE model achieves a relative preference score of 3.56%, indicating a slight preference for its generated paraphrases in terms of meaning preservation. The VAE model performs the best, with a score of 23.00%, while the Separator and Latent BoW models receive negative scores of -17.89% and -8.67%, respectively.
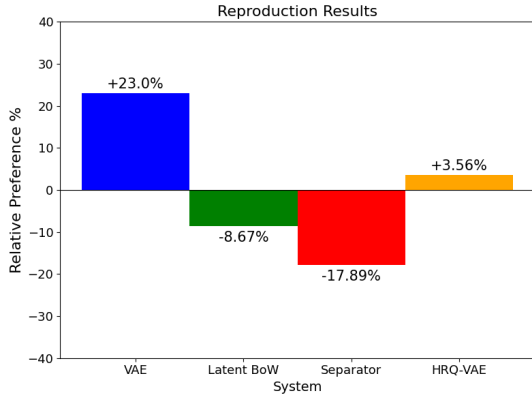
---

Figure 1: Relative preference results from our reproduction.

To facilitate a direct comparison with the original study, we present the results obtained by Hosking et al. (2022) in Figure 2. The original study reports relative preference scores of +36% for the VAE model, -16% for Latent BoW, -24% for Separator, and +4% for HRQ-VAE.
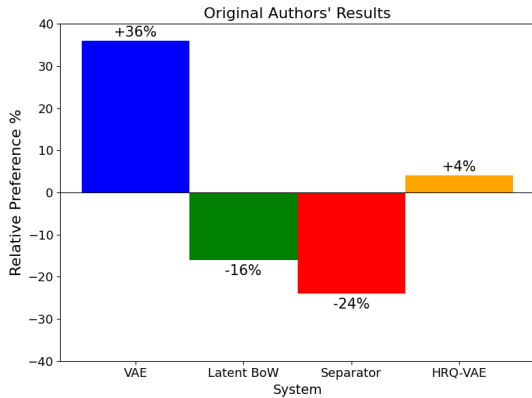


Figure 2: Results obtained by the original authors in Hosking et al. (2022), visualised in a manner consistent with our own findings. The numerical values presented are directly sourced from the original authors' publication.

Comparing the results of our reproduction study with the original findings, we observe some notable differences. While the VAE model maintains its position as the best-performing model in both studies (for preserving meaning), the relative preference scores for the other models vary. In our reproduction, the HRQ-VAE model is slightly less preferred (3.56%) than in the original study (4%). The Separator model is more preferred in our study (-17.89%) compared to the original (-24%), while the Latent BoW model is less preferred (-8.67%) than in the original (-16%), with negative scores indicating less preference for the model. Overall, our replication study shows a narrowing in the range of model

preferences: the best models are not as strongly preferred, and the least preferred models are not as strongly disliked as in the original study, even though the ranking order remains the same.

### 4.1. Quantified Reproducibility Assessments (QRA)

To further evaluate the reproducibility of the original study, we conducted Quantified Reproducibility Assessments (QRA) as described by Belz et al. (2021). These assessments provide a standardised way to quantify the degree of reproducibility between the original study and our reproduction. The code used to do these calculations can be found alongside the data on our github repo[7].

#### 4.1.1. Type I Assessment

Type I assessment measures the reproducibility of individual results using the coefficient of variation (CV*). CV* is an adjusted version of the coefficient of variation that accounts for small sample sizes (Belz, 2022). It can be used even with pairs of results, such as those obtained from an original study and its reproduction. We calculated the CV* for each model by comparing the original and reproduction percentage scores.

$$CV^* = \left(1 + \frac{1}{4n}\right)\frac{s^*}{|\bar{x}|} \qquad (1)$$

where $s^*$ is the unbiased sample standard deviation, $\bar{x}$ is the sample mean, and $n$ is the sample size.

Table 1: Type I (CV*) Assessment

| System | Original | Reproduction | CV* |
|---|---|---|---|
| VAE | +36% | +23% | 43.936 |
| Latent BoW | -16% | -8.67% | 59.246 |
| Separator | -24% | -17.89% | 29.084 |
| HRQ-VAE | +4% | +3.56% | 11.605 |

#### 4.1.2. Type II Assessment

Type II assessment evaluates the reproducibility of a set of results using correlation measures. We calculated Pearson's $r$ and Spearman's $\rho$ correlations between the original and reproduction percentage scores.

Table 2: Type II (Correlation) Assessment

| Metric | Value | *p*-value |
|---|---|---|
| Pearson's $r$ | 0.995 | 0.0049 |
| Spearman's $\rho$ | 1.000 | <0.0001 |

---

[7]https://github.com/NapierNLP/repronlp_2024

## 5. Discussion

Our reproduction study aimed to assess the reproducibility of the human evaluation results reported by Hosking et al. (2022) for their proposed hierarchical sketch induction approach to paraphrase generation. By closely following their methodology but using the Prolific platform for participant recruitment, only screening participants based on location and focusing on the meaning preservation criterion, we sought to determine to what extent the original findings could be replicated.

The results of our reproduction study show a similar trend to the original findings, with the VAE model clearly achieving the highest relative preference score for meaning preservation. However, we observed some notable differences in the magnitudes of the relative preference scores for the other models. The HRQ-VAE model, which was the main focus of the original study, received a slightly lower preference score in our reproduction (3.56%) compared to the original (4%). Additionally, the Separator and Latent BoW models exhibited different degrees of dislike compared to the original study. The Separator model was less disliked in our reproduction, with a relative preference score of -17.89% compared to -24% in the original study. Similarly, the Latent BoW model was also less disliked in our reproduction, receiving a score of -8.67% compared to -16% in the original study.

To further evaluate the reproducibility of the original study, we conducted Quantified Reproducibility Assessments (QRA) as described by Belz et al. (2021). The assessment of individual model reproducibility using the coefficient of variation (CV*) revealed some variability, with the Separator model showing the best reproducibility (CV* = 29.0843) and the Latent BoW model having the lowest reproducibility (CV* = 59.2464). However, the assessment of the overall reproducibility using correlation measures demonstrated a strong positive correlation between the original and reproduction results. Both Pearson's $r$ (0.995, $p = 0.0049$) and Spearman's $\rho$ (1.000, $p < 0.0001$) indicated a high degree of overall reproducibility.

Despite these differences, the overall ranking of the models in terms of meaning preservation remained consistent between the original study and our reproduction. This suggests that the fundamental findings of the original study are reproducible to some extent, even with the modifications made to the participant recruitment platform, and the focus on a single evaluation criterion.

It is important to acknowledge the limitations of our reproduction study. First, the use of a different participant recruitment platform (Prolific) and the exclusion of certain participant screening criteria may have introduced variability in the evaluator pool, potentially influencing the results. Second, focusing on a single evaluation criterion (meaning preservation) rather than the three criteria used in the original study may have simplified the task for participants but also limited the scope of the reproducibility assessment.

## 6. Conclusion

Our findings contribute to the broader discussion on the reproducibility of human evaluation studies in NLP research. The fact that we were able to largely reproduce the original results, despite the modifications made, highlights the potential for reproducing human evaluation findings across different platforms and with variations in the evaluation setup. However, the observed differences in the relative preference scores underscore the sensitivity of human evaluations to factors such as participant recruitment and the specific evaluation criteria used.

To further enhance the reproducibility of human evaluation studies, we recommend that researchers provide detailed documentation of their methodology, including participant recruitment procedures, evaluation guidelines, and analysis methodologies. Additionally, we strongly suggest publishing both raw data and analysis code where possible. This transparency will facilitate replication attempts and enable more robust comparisons across studies. Additionally, exploring the impact of different participant pools and evaluation setups on the reproducibility of results can provide valuable insights into the generalisability of human evaluation findings.

In conclusion, our reproduction study demonstrates that the human evaluation results reported by Hosking et al. (2022) are partially reproducible when using a different participant recruitment platform and focusing on a single evaluation criterion. While we observed some differences in the relative preference scores, the overall ranking of the models remained consistent with the original findings. This study contributes to the ongoing efforts to assess and improve the reproducibility of human evaluation studies in NLP research, and highlights the importance of detailed documentation and exploration of factors influencing reproducibility. Future work should continue to investigate the robustness of human evaluation findings across different setups and participant pools to strengthen the reliability and generalisability of NLP evaluation practices.

## 7. Bibliographical References

Mohammad Arvan and Natalie Parde. 2023. Human evaluation reproduction report for data-to-text generation with macro planning. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 89–96, Varna, Bulgaria.

Anya Belz. 2022. A metrological perspective on reproducibility in nlp. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Shubham Agarwal, Anastasia Shimorina, and Ehud Reiter. 2020. ReproGen: Proposal for a shared task on reproducibility of human evaluations in NLG. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 232–236, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz, Shubham Agarwal, Anastasia Shimorina, and Ehud Reiter. 2021. A systematic review of reproducibility research in natural language processing. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 381–393, Online. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2023. The 2023 ReproNLP shared task on reproducibility of evaluations in NLP: Overview and results. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 35–48, Varna, Bulgaria.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *Proceedings of the Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew Dai, Rafal Jozefowicz, and Samy Bengio. 2016. Generating sentences from a continuous space. In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning*, pages 10–21, Berlin, Germany. Association for Computational Linguistics.

Zihang Chen, Hongbo Zhang, Xiaoji Zhang, and Leqi Zhao. 2017. Quora question pairs.

Anthony Fader, Luke Zettlemoyer, and Oren Etzioni. 2013. Paraphrase-driven learning for open question answering. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1608–1618, Sofia, Bulgaria. Association for Computational Linguistics.

Yao Fu, Yansong Feng, and John P Cunningham. 2019. Paraphrase generation with latent bag of words. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Mingqi Gao, Jie Ruan, and Xiaojun Wan. 2023. A reproduction study of the human evaluation of role-oriented dialogue summarization models. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 124–129, Varna, Bulgaria.

Sebastian Gehrmann, Elizabeth Clark, and Thibault Sellam. 2023. Repairing the cracked foundation: A survey of obstacles in evaluation practices for generated text. *Journal of Artificial Intelligence Research*, 77:103–166.

Javier González Corbelle, Jose Alonso, and Alberto Bugarín-Diz. 2023. Some lessons learned reproducing human evaluation of a data-to-text system. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 49–68, Varna, Bulgaria.

Tom Hosking and Mirella Lapata. 2021. Factorising meaning and form for intent-preserving paraphrasing. 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 1405–1418, Online. Association for Computational Linguistics.

Tom Hosking, Hao Tang, and Mirella Lapata. 2022. Hierarchical sketch induction for paraphrase generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2489–2501,

Dublin, Ireland. Association for Computational Linguistics.

David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 169–182, Dublin, Ireland. Association for Computational Linguistics.

Manuela Hürlimann and Mark Cieliebak. 2023. Reproducing a comparative evaluation of German text-to-speech systems. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 136–144, Varna, Bulgaria.

Takumi Ito, Qixiang Fang, Pablo Mosteiro, Albert Gatt, and Kees van Deemter. 2023. Challenges in reproducing human evaluation results for role-oriented dialogue summarization. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 97–123, Varna, Bulgaria.

Filip Klubička and John D. Kelleher. 2023. HumEval'23 reproduction report for paper 0040: Human evaluation of automatically detected over- and undertranslations. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 153–189, Varna, Bulgaria.

Yiru Li, Huiyuan Lai, Antonio Toral, and Malvina Nissim. 2023. Same trends, different answers: Insights from a replication study of human plausibility judgments on narrative continuations. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 190–203, Varna, Bulgaria.

Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312.

Saad Mahamood. 2023. Reproduction of human evaluations in: "it's not rocket science: Interpreting figurative language in narratives". In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 204–209, Varna, Bulgaria.

Margot Mieskes and Jacob Georg Benz. 2023. h_da@ReproHumn – reproduction of human evaluation and technical pipeline. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 130–135, Varna, Bulgaria.

Ondrej Platek, Mateusz Lango, and Ondrej Dusek. 2023. With a little help from the authors: Reproducing human evaluation of an MT error detector. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 145–152, Varna, Bulgaria.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Craig Thomson, Ehud Reiter, and Anya Belz. 2024. Common Flaws in Running Human Evaluation Experiments in NLP. *Computational Linguistics*, pages 1–11.

Emiel van Miltenburg, Anouck Braggaar, Nadine Braun, Debby Damen, Martijn Goudbeek, Chris van der Lee, Frédéric Tomas, and Emiel Krahmer. 2023. How reproducible is best-worst scaling for human evaluation? a reproduction of 'data-to-text generation with macro planning'. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 75–88, Varna, Bulgaria.

Lewis Watson and Dimitra Gkatzia. 2023. Unveiling NLG human-evaluation reproducibility: Lessons learned and key insights from participating in the ReproNLP challenge. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 69–74, Varna, Bulgaria.

# ReproHum #0043-4: Evaluating Summarization Models: Investigating the Impact of Education and Language Proficiency on Reproducibility

**Mateusz Lango, Patrícia Schmidtová, Simone Balloccu, Ondřej Dušek**

Charles University, Faculty of Mathematics and Physics

Institute of Formal and Applied Linguistics

Prague, Czech Republic

`{lango,schmidtova,balloccu,odusek}@ufal.mff.cuni.cz`

## Abstract

In this paper, we describe several reproductions of a human evaluation experiment measuring the quality of automatic dialogue summarization (Feng et al., 2021). We investigate the impact of the annotators' highest level of education, field of study, and native language on the evaluation of the informativeness of the summary. We find that the evaluation is relatively consistent regardless of these factors, but the biggest impact seems to be a prior specific background in natural language processing (as opposed to, e.g. a background in computer science). We also find that the experiment setup (asking for single vs. multiple criteria) may have an impact on the results.

**Keywords:** human evaluation, reproduction, reproducibility, dialogue summarization, summarization

## 1. Introduction

Human evaluation is generally considered to be the gold standard for Natural Language Processing (NLP) systems assessment. However, many factors can affect its reliability. Subjectivity in human ratings can make experiments impossible to reproduce (Belz et al., 2021); the adopted definition for the evaluated criteria can confuse the annotators (Hosking et al., 2024), and external factors (e.g. fluency) can influence them (Wu et al., 2023). As researchers, we often do not realize the flaws in our own evaluation schemes (Thomson et al., 2024), but they can be found when someone else tries to reproduce such evaluation. Therefore, efforts such as the ReproHum project (Belz et al., 2023a; Belz and Thomson, 2023) also help us design better and more robust human evaluation practices.

In this paper, we describe our attempt at reproducing the human evaluation experiment on dialogue summarization originally presented by Feng et al. (2021) (see Section 2). We specifically target reproduction on one of the datasets and focus mainly on the informativeness criterion. We set up an initial experiment with a setting as close as possible to the original study, including hiring the participants. We then run three variants of the reproduction, inspecting the effect of annotators native language and general background (including knowledge of NLP), using participants hired over the Prolific crowdsourcing platform (see Section 3).

Our reproductions were able to confirm some of the original paper's high-level conclusions from the human evaluation, but also showed some substantial differences in relative rankings among Feng et al. (2021)'s own systems as well as in absolute ratings of all evaluated summarizers (see Sections 4 and 5). The differences between individual reproductions indicate that while participants' native language and general background are not very relevant, a specific background in natural language processing (NLP) can have an impact. In addition, the particular setup of the evaluation (i.e., checking for single or multiple criteria) seems to affect the results (see Sections 6 and 7).

## 2. Original Experiment

The original paper (Feng et al., 2021) proposes a method on how to leverage DialoGPT (Zhang et al., 2020b) as a dialogue annotator to assist in the task of dialogue summarization. The annotations are added as a pre-processing step prior to the summarization.

The authors test their methods on two datasets: SAMSum (Gliwa et al., 2019) and AMI (Carletta et al., 2006). The performance is evaluated using a combination of automatic metrics – ROUGE (Lin, 2004) and BERTScore (Zhang et al., 2020a) – as well as human evaluation. The authors report that their additional DialoGPT-derived annotations are capable of improving the performance of a pre-trained summarizer – BART (Lewis et al., 2020) and a non-pre-trained summarizer – Pointer-Generator Networks (PGN, See et al., 2017) on both datasets. They also report a new state-of-the-art performance on the SAMSum dataset.

While the paper includes results with multiple external baselines and BART pre-trained model extensions, these are either only used with automatic metrics, or only on the SAMSum dataset. For human evaluation on the AMI dataset, which is rel-

evant for our reproduction, the compared systems were:

- Hierarchical Meeting summarization Network (HMNet, Zhu et al., 2020) – a variant of the encoder-decoder transformer model, specially adapted for modelling dialogues. The network was pre-trained on news summarization data. The results of this model serve as a strong external baseline.

- Vanilla Pointer-Generator Network (PGN, See et al., 2017) is an LSTM-based model that combines standard encoder-decoder architecture with pointer network. No pretraining was applied. This model was used by the authors as a baseline summarizer, which was extended with different annotations proposed by the authors.

- PGN with keyword extraction annotation ($D_{KE}$) – the input to the PGN summarizer is extended with a list of keyword words extracted by analysing the outputs of DialoGPT.

- PGN with redundancy detection annotations ($D_{RD}$) – a special tag is added in front of each utterance in the dialogue that has been detected as redundant.

- PGN with topic segmentation annotations ($D_{TS}$) – a special tag is prepended to each utterance that starts a new topic in the dialogue as detected by DialoGPT.

- PGN with all the above annotations ($D_{ALL}$)– the input to PGN is enhanced with all the additional annotations described above.

Additionally, one dialogue summary written by a human was evaluated for comparison.

Human evaluation is performed on summaries generated for 10 randomly selected dialogues. Four annotators are asked to rate informativeness, conciseness and coverage on a 5-point Likert scale, as well as provide a binary good/bad indication for each summary. More details on the experiment are discussed in Section 3, where we also describe key differences of our reproduction.

## 3. Reproduction Studies

We performed four reproductions of the experiment described above: one according to the ReproHum project guidelines (dubbed ReproHum) and three additional ones (dubbed Repro #1 through #3) to investigate different factors influencing the results of human evaluation. We have tried to follow the original experimental setup as closely as possible, but there are still several differences between the original experiment and our reproductions. We first describe the conditions for the ReproHum study, then detail how the additional studies differ from it.

### 3.1. ReproHum reproduction

The ReproHum reproduction used the following setup as a result of the original study's setup and ReproHum guidelines (Belz and Thomson, 2024):

**Datasets**  The original experiment was performed on both SAMSum (Gliwa et al., 2019), and AMI (Carletta et al., 2006), datasets, but the reproduction was limited to the latter dataset only to limit cost. same 10 dialogues from AMI were used, presented in the same order.

**Annotation interface**  The annotation interface was slightly different. In the original study, the authors used a simple text file to collect annotations – more specifically, they used a custom script that printed the dialogues on the console and then prompted the user to rate different summaries. We performed the annotations in a Google form, following ReproHum guidelines.

**Evaluated quality factors**  The original annotations include three quality factors: informativeness, conciseness, coverage, and a final binary rating of whether the presented summary is good. Interestingly, the results of the final binary evaluation were not reported in the original study. In our reproduction, we limited the study to the evaluation of informativeness only, following the decision of the ReproHum team.

**Annotators**  All the original annotators were Chinese PhD students with a background in NLP, specifically in text generation or summarization tasks. Their level of English was assessed by a Chinese state examination of English proficiency: College English Test (CET-6).

In our ReproHum reproduction, we aimed at getting a close demographic, with main focus on hiring PhD students. Our annotators were thus all PhD students and non-native speakers of English, hired on a contract basis. However, they had no background in computer science, their native languages did not include Chinese, and their fluency level was self-assessed.

**Remuneration**  In the original study, the human evaluators were paid 10 USD each. According to the ReproHum fair pay policy, the reproduction wage was set at 14.3 USD per hour. The time needed to perform the annotation was estimated

to be 8 hours, which resulted in a total wage of approx. 115 USD per annotator.[1]

**Annotation guidelines** The original study used annotation guidelines in Chinese. As the annotators in our reproduction did not speak Chinese, we translated the annotation guidelines into English. In addition, since our reproduction concerns only one quality factor and one dataset, the guideline was edited to remove mentions of other quality factors and the SAMSum dataset. The final annotation guideline is as follows:

> *Hi everyone, thanks for helping to do the human evaluation, there is one dataset, AMI, long conversation, 10 data items in total. Please mark each based on the indicator: Informativeness, ranging from 1 to 5, 1 being the worst 5 being the best. Informativeness measures whether the abstract contains the key information from the original conversation. Everyone's document is the same, a total of 4 people will evaluate the same data, and we will then calculate the kappa value to measure the consistency.*

### 3.2. Additional reproductions

We conducted three additional reproduction experiments to investigate the influence of annotators (a) having a background in computer science, (b) having English or Chinese as their first language, (c) answering all four questions as in the original experiment, instead of just one as in the ReproHum reproduction. We followed the same approach as in the ReproHum study, except for annotator demographics and the set of questions (Repro #3 only). We used the Prolific platform as an intermediary to easily find annotators with the necessary background.[2]

The specifics of the additional studies are as follows:

- **Repro #1** was conducted by annotators with a background in computer science (at least a bachelor's degree) and Chinese as their native language.

- **Repro #2** was performed by annotators with a background in computer science (at least having completed a bachelor's degree) and English as their native language.

- **Repro #3** was done by annotators with a background in computer science (at least a bachelor's degree) who were native Chinese speakers. The annotators were responding to all the questions from the original study.[3]

Contrary to the original study, our annotators did not have specific background in NLP and were not PhD students. This difference is given by limited annotator availability on the Prolific platform.

As the workload for an annotator was estimated at 8 hours, we decided to divide the study into 10 parts, corresponding to summaries of 10 evaluated dialogues. Each Prolific annotator was required to complete all parts of the study within a two-week period. Each reproduction was carried out with 4 annotators and the same remuneration as for the ReproHum reproduction.

The task of evaluating long dialogue summaries is not ideally designed for platforms such as Prolific. It relies on reading a long text[4] and then answering several questions on a 1-5 scale (or giving a binary response). There are no attention checks and it is rather difficult to design such. For instance, asking questions about dialogue content could inadvertently suggest to annotators that these parts of the dialogue were important and should be included in the summary. Therefore, we used the time spent on the task as a weak indicator of the annotator's careful reading and analysis of the dialogue content. According to Brysbaert (2019), the average adult has a reading speed of 175 to 300 words per minute (wpm), so annotators who completed the annotation of the first dialogue in a time corresponding to a theoretical reading speed of more than 400 wpm were rejected.

## 4. Main Results

The informativeness values obtained in our reproductions together with the results from the original study (Feng et al., 2021) are presented in Table 1. A rank analysis of these results (the higher, the better) is shown in Table 2.

**Absolute score differences** All the reproductions are very consistent with the original study in rating the informativeness of the human-written summaries highly, at a very similar level. On the contrary, in our reproduction all automatically generated summaries were rated substantially lower than in the original study. As this is consistent regardless

---

[1]The time estimation was done by a trial annotation of two summaries ran by the reproduction authors. The payments were handled in CZK, we provide conversions based on the exchange course as of March 2023 (1 USD = 23.4 CZK).

[2]https://app.prolific.com/

[3]Due to an error in one of the Google Forms, the question about the overall binary quality evaluation was omitted for one dialogue summary.

[4]The joint text of 10 dialogues and the corresponding summaries has almost 65,000 words, which corresponds to 159 A4 pages in 11pt Courier New font.

| | | **Reproductions** | | | |
|---|---|---|---|---|---|
| | Original | ReproHum | Repro #1 | Repro #2 | Repro #3 |
| Evaluated factors | All | Inform. | Inform. | Inform. | All |
| Educational level | PhD Student | PhD Student | $\geq$Bachelor | $\geq$Bachelor | $\geq$Bachelor |
| Background | NLP | Any | CS | CS | CS |
| First language | Chinese | non-English | Chinese | English | Chinese |
| Annotators | In-lab | External | Prolific | Prolific | Prolific |
| Human summary | 4.70 | 4.60 | 4.65 | 4.70 | 4.68 |
| PGN | 2.92 | 1.53 | 1.60 | 1.90 | 1.88 |
| HMNet | **3.52** | **2.68** | **2.23** | **2.90** | **3.08** |
| PGN($D_{KE}$) | 3.20 | <u>1.93</u> | 1.63 | 1.93 | 2.35 |
| PGN($D_{RD}$) | 3.15 | 1.90 | <u>1.75</u> | 1.98 | <u>2.53</u> |
| PGN($D_{TS}$) | 3.05 | 1.85 | 1.60 | 1.98 | 2.38 |
| PGN($D_{ALL}$) | <u>3.33</u> | 1.85 | 1.65 | <u>2.10</u> | 2.18 |
| Fleiss' $\kappa$ | 0.48 | 0.19 | 0.20 | 0.13 | 0.05 |
| Krippendorff's $\alpha$ | | 0.65 | 0.66 | 0.58 | 0.38 |

Table 1: The average informativeness values obtained in the original study and performed reproductions.

| | **Ranks of the final results** | | | | | **Averaged ranks** | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Original | ReproHum | R#1 | R#2 | R#3 | ReproHum | R#1 | R#2 | R#3 |
| Human summary | 7 | 7 | 7 | 7 | 7 | 6.84 | 6.81 | 6.90 | 6.60 |
| PGN | 1 | 1 | 1.5 | 1 | 1 | 2.59 | 3.20 | 3.00 | 2.70 |
| HMNet | **6** | **6** | **6** | **6** | **6** | **5.10** | **4.56** | **5.16** | **4.84** |
| PGN($D_{KE}$) | 4 | <u>5</u> | 3 | 2 | 3 | <u>3.53</u> | 3.43 | 3.10 | 3.43 |
| PGN($D_{RD}$) | 3 | 4 | <u>5</u> | 3.5 | <u>5</u> | 3.48 | <u>3.50</u> | 3.16 | <u>3.81</u> |
| PGN($D_{TS}$) | 2 | 2.5 | 1.5 | 3.5 | 4 | 3.28 | 3.21 | 3.20 | 3.49 |
| PGN($D_{ALL}$) | <u>5</u> | 2.5 | 4 | <u>5</u> | 2 | 3.20 | 3.29 | <u>3.48</u> | 3.14 |

Table 2: The ranked results of informativeness (the higher, the better) obtained in the original study (Original) and performed reproductions (ReproHum, R#1-R#3). We report both the ranks of the averaged informativeness from Table 1, as well as the ranks of informativeness averaged over all samples.

of whether the annotators have a background in computer science, are native English speakers or have a higher level of education, it seems that the main factor influencing this result is the participants' background in NLP (or in NLG tasks in particular) and potential prior experience with automatic summarizers.[5]

**System ranking**  Among the automatically generated summaries, HMNet is consistently assessed as the best method for producing informative summaries and PGN as the worst. The PGN extensions are almost always all ranked in between the basic PGN and HMNet, but their ranking relative to each other varies greatly in different reproductions. This is because there are small absolute differences between them: the standard deviation

of different PGN extensions' results is $\leq 0.14$ for both the original study and all reproductions (even lower, at 0.04-0.07 for ReproHum, Repro #1 and Repro #2 reproductions).

**Inter-annotator agreement**  The inter-annotator agreement is much lower in the reproduced studies as compared to the original experiment – Fleiss' $\kappa$ is in the 0.1-0.2 range instead of the original 0.48. For Repro #3, $\kappa$ is even lower. After looking at the correlation matrix between different annotators, we discovered that the responses of one annotator were poorly correlated with all the other annotators. We investigated the time spent on the annotation, but it was not different from the other annotators. The annotator also ranked human written summaries relatively higher than the other assessed summaries. Nevertheless, after excluding this annotator[6] the Fleiss' $\kappa$ went up to 0.16, taking a value similar to that obtained in other reproductions.

---

[5]The original study was conducted in 2021, before the popularity of ChatGPT, which can also serve as a summarization engine. This may have raised annotators' expectations of the output quality of an AI-based system. However, we have no information on whether our annotators ever used ChatGPT for summarization.

[6]Repro #3 results recomputed for 3 annotators only are presented in Table 3).

**Comparing different reproductions** We do not observe very large differences that would indicate a significant impact of the factors influencing the selection of annotators. The reproduction that seems to stand out the most is Repro #3 (evaluation of all quality factors). As already mentioned, the informativeness of the PGN variants shows larger rating differences in this case (even when excluding the poorly correlated annotator).The absolute rating values are also consistently higher for all the methods, closer to the reproduced study. This may indicate that annotators responding to multiple quality criteria are more likely to try to split the overall quality rating into multiple factors than when presented with a single quality question. However, the observed differences against any other reproduction are not statistically significant.

**Statistical analysis** We performed a statistical analysis of the obtained informativeness results in all reproductions. Following the recommendations of Demšar (2006), we performed the non-parametric Friedman rank test with Nemenyi post-hoc analysis. For all reproductions, the null hypotheses of Friedman tests about the lack of differences in informativeness among all investigated summaries were rejected with low p-values ($p < 0.001$ for all reproductions). The results of the post-hoc analysis are presented in Figure 1 in the form of critical distance plots.

In all reproductions, the differences between the PGN baseline and all the PGN variants with additional annotations proposed by Feng et al. (2021) were not statistically significant at the $\alpha = 5\%$ significance level. In fact, the difference between human-written summaries and summaries provided by HMNet, the best automatic method, was also not significant due to the small sample size. In the main reproduction (ReproHum) and Repro #2 and #3, there is a statistically significant difference between HMNet and the PGN baseline. In contrast, the differences between HMNet and the enhanced variants of PGN are not statistically significant (except for PGN($D_{KE}$) in Repro #2). In Repro #1, all automatic summarization methods are statistically indistinguishable.

**Additional results from Repro #3** The results of Repro #3 include not only the informativeness values, but also the measurements of conciseness, coverage and the assessment of overall evaluation (the latter not being reported in the original work). The results are presented in Table 3. As mentioned earlier, the responses of one of the annotators were poorly correlated with those of the other three annotators, therefore we report the results averaged over all annotators (R#3) and the results averaged over three annotators only (R#3*). The discussion

of the results will focus on the latter variant.

The general observation that our annotators evaluated all systems lower than in the original study remains true for coverage, but we obtained values of similar magnitude for conciseness. The ranking of the best performing methods resulting from the reproduction is similar to the original one for informativeness and coverage (Spearman correlations of 0.75 and 0.79, respectively) but differences are visible for conciseness (Spearman 0.39). Inter-annotator agreement is significantly lower than in the original study for all measures.

Looking at the overall binary quality evaluation, it seems that the PGN baseline is very weak, as none of the produced summaries were rated as good. The extensions of PGN improve the performance, but still fall significantly behind HMNet. Analysing the results of all measures, it seems that $D_{RD}$ is the main cause of the improvement and combining it with other techniques ($D_{ALL}$) does not lead to further improvements, but, on the contrary, degrades the summaries.

## 5. Quantifying Reproducibility

Following the guidelines of the ReproHum shared task (Belz et al., 2023b, Sect. A5), we identify reproduction targets in the following categories:
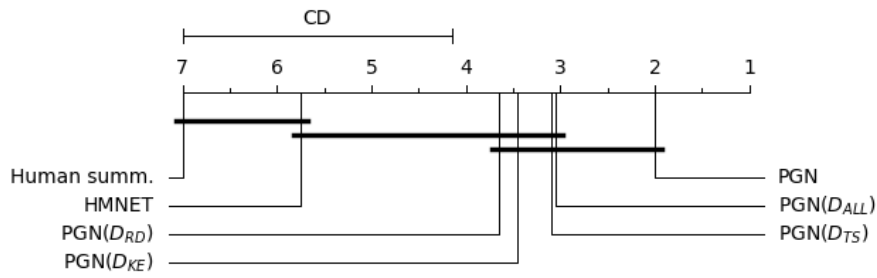
- Type I – numerical scores: the average informativeness of summaries generated by different methods

- Type II – sets of numerical values: the set of informativeness results for all the methods in the study

**Type I** Following the quantified reproducibility assessment by Belz et al. (2022), we computed the small sample coefficient of variation (CV*) as a measure of the degree of reproducibility for numerical scores. The results are given in Table 4.
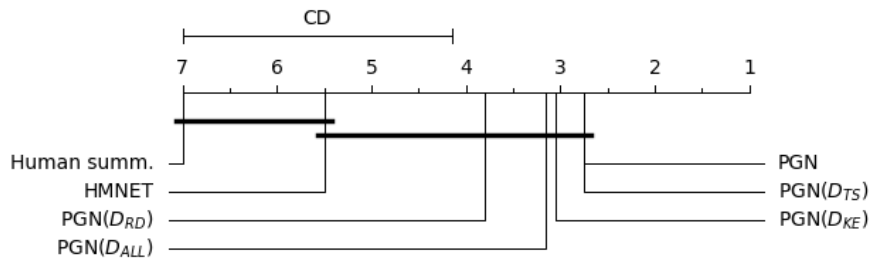
The values of CV* computed for the original study and the main ReproHum reproduction are in the range of 48-63, except for the significantly lower values for HMNet and the summaries written by humans.

As to be expected, the coefficients of variation are smaller when computed for all the performed reproductions and the original study. Most CV* values are in the range of 28-33, again with the exceptions for HMNet and human summaries.
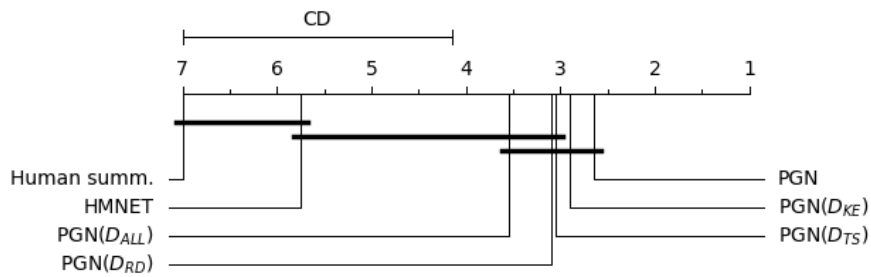
**Type II** results are evaluated with Pearson and Spearman correlation (Huidrom et al., 2022), as well as with the root-mean-square deviations from the original results. The results are presented in Table 5.
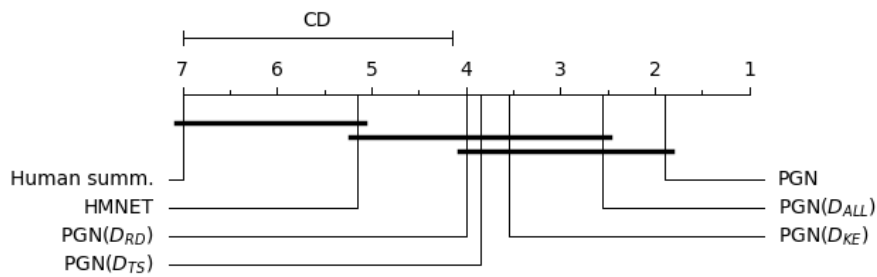
(a) ReproHum

(b) Reproduction 1

(c) Reproduction 2

(d) Reproduction 3

Figure 1: Critical distance diagrams showing the results of post-hoc Nemenyi tests performed for informativeness values obtained in the four performed reproductions. For all reproductions, the global Friedman test rejected the null hypothesis with $p < 0.001$. Critical distance plots present the average rank obtained in the Friedman test (the higher, the better) and show the difference between ranks that would imply statistical significance in the post-hoc analysis (critical distance). If the difference between the methods is not statistically significant, their results are connected with a thick horizontal line. More details on these plots can be found in (Demšar, 2006).

| | Informativeness | | | Conciseness | | | Coverage | | | Overall | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Orig. | R#3 | R#3* | Orig. | R#3 | R#3* | Orig. | R#3 | R#3* | R#3 | R#3* |
| Human summary | 4.70 | 4.68 | 4.97 | 3.85 | 4.28 | 4.43 | 4.35 | 4.58 | 4.80 | 0.98 | 1.00 |
| PGN | 2.92 | 1.88 | 1.53 | 3.08 | 2.55 | 2.27 | 2.70 | 2.05 | 1.67 | 0.13 | 0.00 |
| HMNet | **3.52** | **3.08** | **2.80** | 2.40 | 3.00 | 2.97 | **3.40** | **3.18** | **3.00** | **0.50** | **0.40** |
| PGN($D_{KE}$) | 3.20 | 2.35 | 2.27 | 3.08 | **3.23** | **3.10** | 3.00 | 2.33 | 2.10 | 0.13 | 0.07 |
| PGN($D_{RD}$) | 3.15 | <u>2.53</u> | <u>2.53</u> | **3.25** | <u>3.18</u> | **3.10** | 3.00 | <u>2.53</u> | <u>2.53</u> | 0.13 | 0.10 |
| PGN($D_{TS}$) | 3.05 | 2.38 | 2.17 | 3.10 | 3.03 | 2.87 | <u>3.17</u> | 2.33 | 2.13 | 0.11 | 0.04 |
| PGN($D_{ALL}$) | <u>3.33</u> | 2.18 | 1.90 | **3.25** | 2.85 | 2.70 | 3.10 | 2.08 | 1.80 | 0.10 | 0.10 |
| Fleiss' $\kappa$ | 0.48 | 0.05 | 0.16 | 0.40 | 0.01 | 0.03 | 0.41 | 0.03 | 0.11 | 0.47 | 0.61 |
| Krippendorff's $\alpha$ | | 0.38 | 0.51 | | 0.13 | 0.15 | | 0.35 | 0.45 | 0.47 | 0.61 |

Table 3: The average informativeness, conciseness, coverage and overall binary evaluation of summaries as obtained in the original human evaluation (Orig.) and our Repro #3 (R#3). Additionally, we also report reproduction results computed on data from 3 annotators only (R#3*) - see more details in the text.

| CV* | ReproHum | All Repro. |
|---|---|---|
| Human summary | 2.14 | 1.01 |
| PGN | 62.28 | 31.71 |
| HMNet | 27.02 | 18.51 |
| PGN($D_{KE}$) | 49.36 | 30.91 |
| PGN($D_{RD}$) | 49.36 | 28.51 |
| PGN($D_{TS}$) | 48.83 | 29.11 |
| PGN($D_{ALL}$) | 56.97 | 32.86 |

Table 4: The small-sample coefficient of variation (CV*) of informativeness computed for original and ReproHum study (2 samples) and for all the reproductions (5 samples).

| | Pearson | Spearman | RMSE |
|---|---|---|---|
| ReproHum | 0.99 | 0.85 | 1.16 |
| Repro #1 | 0.98 | 0.88 | 1.35 |
| Repro #2 | 0.98 | 0.88 | 1.00 |
| Repro #3 | 0.97 | 0.68 | 0.77 |

Table 5: The values of root-mean-square deviation, Pearson and Spearman correlations computed between the original and reproduced results.

The Pearson correlations are very high for all the reproduction studies, which can be attributed to the fact that the human summary scores are relatively high outliers in all the studies (after removing them, the correlations drop from 0.97-0.99 to 0.78-0.88). This is also reflected in the lower Spearman correlations, which are more robust to outliers. The lowest Spearman correlation was obtained for Reproduction 3 (0.68) which is the only correlation in this study that is not statistically significant ($\alpha = 5\%$). Note that the sample size is very small (7).

Finally, RMSE values of around 1 reflect the general tendency of our annotators to rate automatic summaries lower than in the original study. The closest results to the original study, as measures by RMSE, were obtained in the Reproduction 3 where all quality factors were evaluated.

## 6. Summary

From the results of the original study, the authors draw three major conclusions (see Sec. 4.5 in Feng et al., 2021):

1. "HMNet gets the best score in informativeness and coverage", which was confirmed by our reproductions.

2. "Our method can achieve higher scores in all three metrics", which again is in line with the results of our reproductions.

3. "We also find there is still a gap between the scores of generated summaries and the scores of golden summaries" – which was not only confirmed in our reproductions, but also the gap seems substantially larger than in the original study.

Nevertheless, the results of the original study also provided evidence that the combination of all proposed annotations ($D_{ALL}$) gives the best informativeness among the PGN variants and that the gap against the better performing HMNet is relatively small (0.19). This was not confirmed by our reproductions. $D_{ALL}$ was the worst PGN extension evaluated in two reproductions, and the best and second best in the other two reproductions. Similarly, the reported gap between the best PGN extension and HMNet ranged from 0.48 to 0.8 on a 5-point scale, at least two and a half times larger than in the original study.

# 7. Discussion

We can attempt to draw some conclusions from the analysis of the differences between our reproductions: Firstly, mother tongue, level of education or field of study do not seem to have a significant impact on the results of human evaluation in the summarisation task. The only exception is a very specific background in NLP technologies. Second, when working on reproduction experiments, it might be better to evaluate all quality factors, even if were are interested in reproducing the result for a single quality factor in particular. Finally, we believe that it is always helpful to carry out a statistical analysis of the results obtained. Even if the analysis is not conclusive, e.g. due to the small sample size, it gives a much better picture of the variability of the results and the conclusions that can be drawn from them.

# 8. Acknowledgements

# 9. Bibliographical References

Anya Belz, Maja Popovic, and Simon Mille. 2022. Quantified reproducibility assessment of NLP results. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16–28, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz, Anastasia Shimorina, Shubham Agarwal, and Ehud Reiter. 2021. The ReproGen shared task on reproducibility of human evaluations in NLG: Overview and results. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 249–258, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2023. The 2023 ReproNLP shared task on reproducibility of evaluations in NLP: Overview and results. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 35–48, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023a. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *Proceedings of the Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Jackie Cheung, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Huiyuan Lai, Chris van der Lee, Emiel van Miltenburg, Yiru Li, Saad Mahamood, Margot Mieskes, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Pablo Mosteiro Romero, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023b. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in nlp.

Marc Brysbaert. 2019. How many words do we read per minute? a review and meta-analysis of reading rate.

Janez Demšar. 2006. Statistical comparisons of classifiers over multiple data sets. *J. Mach. Learn. Res.*, 7:1–30.

Xiachong Feng, Xiaocheng Feng, Libo Qin, Bing Qin, and Ting Liu. 2021. Language model as an annotator: Exploring DialoGPT for dialogue summarization. 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 1479–1491, Online. Association for Computational Linguistics.

Tom Hosking, Phil Blunsom, and Max Bartolo. 2024. Human feedback is not gold standard.

Rudali Huidrom, Ondřej Dušek, Zdeněk Kasner, Thiago Castro Ferreira, and Anya Belz. 2022. Two reproductions of a human-assessed comparative evaluation of a semantic error detection system. In *Proceedings of the 15th International Conference on Natural Language Generation: Generation Challenges*, pages 52–61, Waterville, Maine, USA and virtual meeting. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Craig Thomson, Ehud Reiter, and Anya Belz. 2024. Common Flaws in Running Human Evaluation Experiments in NLP. *Computational Linguistics*, pages 1–11.

Zixiu Wu, Simone Balloccu, Ehud Reiter, Rim Helaoui, Diego Reforgiato Recupero, and Daniele Riboni. 2023. Are experts needed? on human evaluation of counselling reflection generation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6906–6930, Toronto, Canada. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020a. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020b. DIALOGPT : Large-scale generative pre-training for conversational response generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278, Online. Association for Computational Linguistics.

Chenguang Zhu, Ruochen Xu, Michael Zeng, and Xuedong Huang. 2020. A hierarchical network for abstractive meeting summarization with cross-domain pretraining. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 194–203, Online. Association for Computational Linguistics.

## 10.  Language Resource References

Carletta, Jean and Ashby, Simone and Bourban, Sebastien and Flynn, Mike and Guillemot, Mael and Hain, Thomas and Kadlec, Jaroslav and Karaiskos, Vasilis and Kraaij, Wessel and Kronenthal, Melissa and Lathoud, Guillaume and Lincoln, Mike and Lisowska, Agnes and McCowan, Iain and Post, Wilfried and Reidsma, Dennis and Wellner, Pierre. 2006. *The AMI Meeting Corpus: A Pre-announcement*. Springer Berlin Heidelberg.

Gliwa, Bogdan and Mochol, Iwona and Biesek, Maciej and Wawer, Aleksander. 2019. *SAMSum Corpus: A Human-annotated Dialogue Dataset for Abstractive Summarization*. Association for Computational Linguistics.

## A.  Human Evaluation Datasheet (HEDS)

Human Evaluation Datasheet (HEDS, Shimorina and Belz, 2022) for the main ReproHum reproduction (see Sec. 3.1) is provided in the ReproHum GitHub repository.[7]

---

[7]https://github.com/nlp-heds/repronlp2024

# ReproHum #0033-3: Comparable Relative Results with Lower Absolute Values in a Reproduction Study

## Yiru Li, Huiyuan Lai, Antonio Toral, Malvina Nissim

CLCG, University of Groningen
Groningen, the Netherlands
y.li.170@student.rug.nl
{h.lai, a.toral.ruiz, m.nissim}@rug.nl

## Abstract

In the context of the ReproHum project aimed at assessing the reliability of human evaluation, we replicated the human evaluation conducted in "*Generating Scientific Definitions with Controllable Complexity*" by August et al. (2022). Specifically, humans were asked to assess the fluency of automatically generated scientific definitions by three different models, with output complexity varying according to target audience. Evaluation conditions were kept as close as possible to the original study, except of necessary and minor adjustments. Our results, despite yielding lower absolute performance, show that relative performance across the three tested systems remains comparable to what was observed in the original paper. On the basis of lower inter-annotator agreement and feedback received from annotators in our experiment, we also observe that the ambiguity of the concept being evaluated may play a substantial role in human assessment.[*]

**Keywords:** human evaluation, reproducibility, ReproHum

## 1. Introduction

In spite of substantial advances in the automatic evaluation of Natural Language Processing (NLP), especially with the development of trained metrics highly correlating with human judgements, such as COMET (Rei et al., 2020), eventually it is the actual human evaluations that are still widely considered the most significant and reliable performance assessment. This is particularly true in language generation tasks, where the availability of a human gold standard produced in advance, as it is common practice in classification tasks, is not an option due to the large variability of valid outputs.

And yet, human evaluation, both in classification and generation tasks, is surely not free of problems. First, humans might not be great judges on a given task as they cannot tell one category from another; this has been shown for example in profiling (Flekova et al., 2016), in the detection of political leaning (De Mattei et al., 2020), and in discerning AI-generated from human-written texts (Clark et al., 2021; Freitag et al., 2021). Second, even when people might be able to yield judgements in a given task, how to perform human evaluations which are dependable, for example on conversations (Smith et al., 2022), is an open problem. Third, and most importantly, human judgements are tainted by a somewhat natural variability, which might yield idiosyncratic results that are not reproducible in subsequent studies and thus eventually not that indicative of system performance beyond a specific and single experiment. This is especially true if evaluation settings are not systematically and clearly defined and reported. Recent research has shown that due to these and related factors, reproducing human evaluation in NLP studies proves an almost impossible task (Belz et al., 2023).

This paper situates itself in this last line of research, in the context of the larger ReproHum[1] project (Belz and Thomson, 2024), which is a multi-lab cooperative project aiming to test the reproducibility of human evaluations through large-scale reproductions.

Our reproduction work follows the schedule provided by the project coordination team, and this paper reports our results accordingly. The experiment was pre-registered through the Human Evaluation Data Sheet (HEDS[2]) proposed in Shimorina and Belz (2022)'s work, providing records for possible future usage. In this report, we first summarize the original study and provide a detailed explanation of the human evaluation task we are reproducing (Section 2). Next, we introduce the adjustments we had to make to successfully replicate the experiment (Section 3). Lastly, we report our results and bring forward our observations and comments on the feasibility and meaning of this reproduction (Section 4).

---

[*]In the ReproHum project this reproduction study has code #0033-3.

[1]https://reprohum.github.io/

[2]Details in Appendix A, also on https://github.com/nlp-heds/repronlp2024

## 2. Original Study

The original study we have reproduced is one of the human evaluation tasks described in the paper "*Generating Scientific Definitions with Controllable Complexity*" by August et al. (2022). This research proposes a new method for generating scientific definitions with controllable complexity, varying according to target audience. Several systems are trained using a newly collected dataset of scientific definitions and both automatic and human evaluations are performed on the generated outputs.

### 2.1. Task and Model

The core task in the research is to generate scientific definitions with controllable complexity that are appropriate answers to a "term question" in the form of "What is (are) X," where X is a scientific term or concept (August et al., 2022, Section 3). In the first part of their paper, the authors explore the performance of different models in generating scientific definitions without complexity control. Pairs of the "term questions" and corresponding definitions are then used as training/finetuning data for multiple language models. The authors have also collected additional data from scientific abstracts serving as supporting documents. Through the use of automatic metrics, they conclude that the `BART` model (Lewis et al., 2020) trained with term question concatenated with the supporting document (`BART SD`) outperforms the rest of the models they tested. Therefore, `BART SD` is used as the base generation model for all subsequent experiments.

After the selection of the base generation model, the authors explore four complexity control methods, including their proposed new method called *reranking*. A *Reranker* is composed of two parts: a `BART SD` generator that provides 100 definitions of the same scientific question, and a discriminator that was trained to distinguish scientific journals from science news. The logits of the discriminator are then used to determine the complexity of the definitions. In their work, the original authors have trained one model for each method other than *reranking*, and two models using *reranking* - one of which uses a Linear SVM Classifier as the discriminator and the other one uses the SciBERT uncased pretrained model (Beltagy et al., 2019).

Models representing the four complexity control methods are trained to provide definitions of either high complexity or low complexity and the resulting definitions are then evaluated by means of automatic metrics. See Table 1 for an example of generated definitions, directly taken from August et al. (2022).

### 2.2. Human Evaluation Task

The original paper includes several human evaluation tasks on the generated definitions to test the robustness of their proposed *reranking* approach. 50 terms were randomly selected from the test split as target terms. The corresponding definitions generated for these 50 terms, both with high and low complexity, by the three models that showed the best performance in the automatic evaluation task are then put through human evaluation. These three models are *Reranker* utilizing an SVM classifier as the discriminator, the Generative discriminators (GeDi) proposed by Krause et al. (2021), and the Ensemble of language models (DExperts) proposed by Liu et al. (2021). These $50 \times 2 \times 3 = 300$ definitions were then rated by human annotators.

Besides the main evaluation task that targets the complexity of definitions generated by different systems, three additional side human evaluation tasks are conducted to ensure the generations are all fluent, relevant to the questions, and factual. In our reproduction study we only focused on one of the side evaluation tasks: fluency.

Two trained annotators performed the fluency task, but the specific training they underwent was left unspecified in the original paper. One annotator was one of the authors of the original paper, and the other annotator was a research assistant.

For the evaluation, the annotators were first shown an instruction page, received instructions on the nature of this task. They were informed that they will be evaluating definitions based on their fluency, and the answers to the question "How fluent is this definition?", will need to be given on a scale from "Not at all" to "Very". They were provided with two "Not at all" fluent definitions and two "Very" fluent definitions. Then, on the annotation interface, more specific instructions were given.

In the actual annotation part, the annotators are asked to evaluate to what extent a definition is fluent. They are expected to rate the definitions using 1–4 Likert scales (1 = "Not at all" to 4 = "Very"). The two annotators rate the same 300 definitions. The definitions are presented to them in different, random orders. It is not disclosed to them whether a definition is of high complexity or low complexity, nor which system has generated this definition. In the case of "nan" (empty definition), they should rate it as "Not at all" fluent.

Figure 1 and Figure 2 show the original instruction page and an example of the original annotation interface, respectively. Note that these screenshots were taken from the interface utilized by the main evaluation task on complexity performed in the original experiment, as we do not possess the actual interface used in the fluency evaluation task. Nevertheless, according to the information provided by the original authors, the interfaces used by these

| Control Method | Direction | |
| --- | --- | --- |
| | **Low (News)** | **High (Journal)** |
| SVM-Rerank | A type of computing in which there are many computers running at the same time in different parts of the world. | In computer science, distributed computing is the process of computing on a large scale **without a single centralized data center**. |
| BERT-Rerank | A type of computer system in which there are more than a few computers working together. | In computer science, distributed computing is the process of computing on a large scale **without a single centralized data center**. |
| GeDi | Is the implementation of computer programs across multiple computers on similar hardware and/or software resources. | In computer science, **a concept that states that data must be shared across computing resources**. |
| DExpert | An **Internet-driven by-computing** that portion of different computers from start to finish. | In computer science and communication between-Consequently-integrates. |
| PPLM | **Easeless, self-organized, and often self-organizing** networked computer systems intended for the purposes of optimization. | Multi-purpose, distributed **system software** with or without a single datum storage system. |

Table 1: Generated definitions from each complexity control method for the question: What is (are) distributed computing? Factually incorrect information is labeled in **bolded red**.
*Note*: From "*Generating scientific definitions with controllable complexity*" by August et al. (2022).

two tasks are identical except for the task-specific instructions and questions.

Based on the results, the authors conclude that their SVM-reranked methods can provide definitions that were rated close to "Very" fluent and are significantly more fluent compared to definitions generated by the other two systems. Further discussion of their results, also in comparison with ours, is included later in Section 5.

## 3.   Reproduction Study

In our reproduction study, several adjustments had to be made for various reasons. None of these adjustments are related to the nature of the assessment questions - they remained identical to what was given in the original experiment.

The first adjustment we made was changing the evaluation platform from *LabintheWild* to *Qualtrics*, essentially leading to the re-writing of the evaluating interface. By the time we started reproducing the experiment, *LabintheWild* was inaccessible through its website, forcing us to use another evaluation platform instead; we chose *Qualtrics* since it could replicate the functionality and look-and-feel of the original interface, and we are familiar with it. We tried our best to keep the new interface as

similar as possible to the original interface, keeping important features identical. Figure 3 shows our instructions, and Figure 4 shows an example of our new annotation interface. It is important to note that the instructions for the fluency evaluation task were not reported in the paper nor in the additional information kindly provided through email by the paper's author. As the instruction screenshots provided to us only included examples for the complexity evaluation task, we could not replicate what was included in the original instructions and had to include new examples in our guidelines.

The second adjustment we made was removing other unrelated questions from the interface, now giving our annotators one question per page instead of two questions per page. This change is due to the fact that we are only replicating the fluency evaluation task but not the relevance evaluation task which is included in the original paper alongside the fluency one. The annotators in our replication study are now answering only 300 questions in total (one question – fluency – per instance) instead of 600 in the original paper (two questions – fluency and relevance – per instance). Even though the other 300 questions/answers are irrelevant to the fluency evaluation task, the annotators' overall performance may still be affected by this difference,

# Instructions

You will be given 3 terms with their definitions and asked to rate how complicated and understandable the definitions are.

You will be asked to rate the how complicated and understandable the definition is on a scale from **Not at all** to **Very**.

Examples of very complicated definitions:

**Term:** Acanthoma

**Definition:** An acanthoma is a skin neoplasm composed of squamous or epidermal cells. It is located in the prickle cell layer.

**Term:** Transformer

**Definition:** The Transformer is a deep learning model architecture relying entirely on an attention mechanism to draw global dependencies between input and output.

Examples of not at all complicated definitions:

**Term:** Acanthoma

**Definition:** An acanthoma is a small, reddish bump that usually develops on the skin of an older adult.

**Term:** Transformer

**Definition:** The Transformer is a program used by computers to weigh the importance of different parts of data.

**Please do not press the back button while taking this task.**

Figure 1: A screenshot of the original instruction page.

---

You are currently on section: 1 / 300

### Instructions

**Please read the following text and answer the questions below.**
When rating definitions, please focus on unfamiliar terms or very long, complicated sentences, not grammar.

If a definition's text only says 'nan', please rate it as **Very** complex and **Very** hard to understand.

**Term:** etchplain

**Definition:** See plain.

\* How complicated is the definition's text?

Not at all ○ ○ ○ ○ Very

\* Imagine you are looking up this term, how hard is it for you to understand this definition?

Not at all ○ ○ ○ ○ Very

This includes the definition having terms that are unfamiliar to you.

Figure 2: A screenshot of the original annotation interface.

241

Figure 3: A screenshot of the instruction page in our replication study.



Figure 4: A screenshot of the annotation interface used in our replication study.

however to what extent is unknown.

The third adjustment we made concerned the annotators. We provided the annotators with monetary compensation and they were not related to

this reproduction paper, i.e. none of the annotators is a coauthor. The amount of monetary compensation was determined according to the minimum wage in the U.K. in December 2023. Given the assumption that the annotation task should take approximately 2.5 hours to complete, each annotator was paid 34.6 euros. As said, in the original study, one author of the paper participated in the annotation process; according to the ReproHum reproduction instructions, we have not included one of us in the evaluation task, but instead recruited one NLP PhD student and one linguistics researcher for the task, trying to match as close as possible the background of the original annotators. This adjustment may have had a larger influence on the result than the other modifications we have described: despite the original paper stating that none of their annotators have seen the generations to be evaluated before their evaluation exercise, their familiarity and association with the project could have unintentionally affected the evaluation results.

## 4. Results

### 4.1. Side-by-side Presentations

Table 2 shows a side-by-side presentation of our results and the original results.

|  | **Original** | **Replication** |
| --- | --- | --- |
| Fluency (s.d) SVM-Reranker | 3.71 (0.59) | 3.02 (1.10) |
| Fluency (s.d) GeDi | 3.20 (1.06) | 2.40 (1.20) |
| Fluency (s.d) DExpert | 2.33 (0.85) | 1.81 (1.04) |
| t-test between SVM & GeDi | $t_{198} = 5.99$, $p < 0.001*$, Cohen's $d = 0.60$ | $t_{198} = 4.42$, $p < 0.001*$, Cohen's $d = 0.62$ |
| t-test between SVM & DExpert | $t_{198} = 18.85$, $p < 0.001*$, Cohen's $d = 1.88$ | $t_{198} = 9.65$, $p < 0.001*$, Cohen's $d = 1.36$ |

Table 2: Comparison of original and reproduction results. * =$p$-value corrected for multiple hypothesis testing using the Bonferroni-Holm correction.

### 4.2. Quantified Reproducibility Assessments

According to the Common Approach of Reproduction provided by the ReproHum Team, we report the three quantified reproducibility assessments below, including adjusted Coefficient of Variation (CV*), Pearson's $r$, and Krippendorff's $\alpha$.

We report an adjusted version of the Coefficient of Variation (CV*) as mentioned in Belz et al. (2022)'s work on quantified reproducibility assessments. CV* was specifically adjusted for small samples. As the experiment utilized a Likert scale from 1-4, we shifted the values from [1,4] to [0,3] to meet the requirement of utilizing CV*. We report the two-way CV* values in Table 3.

| **System** | **CV*** |
| --- | --- |
| SVM-Reranker | 29.09 |
| GeDi | 44.31 |
| DExpert | 48.45 |

Table 3: Two-way CV* between original results and replication results

We have calculated Pearson's correlation coefficient between the original results and the reproduction results as instructed. However, it is worth noting that since the sample size in our case is extremely small ($n = 3$), the coefficient (Pearson's $r = .987$) is not reliable. Spearman's $\rho$ is not suitable for such a small sample size either.

To compare the inter-annotator agreement, we also report the Krippendorff's $\alpha$ of our annotations. The original study reports Krippendorff's $\alpha = 0.63$, while our study reports Krippendorff's $\alpha = 0.45$.

## 5. Discussion and Conclusion

Through the analysis of results, we observe that our results support the finding in the original paper, that the definitions produced by the SVM-reranked method are significantly more fluent compared to definitions generated by the other two systems evaluated. However, we observe that in our reproduction experiment the overall fluency is rated lower for all three systems. In one of last year's ReproHum reports on a different reproduction study (Li et al., 2023), the authors noticed the same phenomenon: The reproduction results support the comparative statements made in the original paper (e.g., one system performs better than the others) with the same overall trend, but with lower overall scores. As the fluency score of SVM-Reranker in our evaluation did not surpass 3.5 as it did in the original experiment, we could not confirm the statement suggested in the original paper that the SVM-Reranker can be rated as nearly "Very fluent".

The two-way CV* values suggest medium to low reproducibility, while the reproduced annotations on definitions generated by SVM-Reranker seem

to have a higher agreement with the original annotations, compared to annotations on other models' definitions. We have also noted a decline in Krippendorff's $\alpha$. The decline of inter-annotator agreement may be attributed to the fact that the original annotators were "trained", while we did not train our annotators since the training process was not specified in the original paper. From the feedback we received from our annotators, the definition of fluency remained ambiguous to some extent, even with the examples and instructions. As a result, the different understanding of the concept of fluency may have caused our two annotators to disagree on a few questions. Lastly, the fact that one of the annotators in the original study was one of the paper's co-author might have influenced the original agreement and thus contributed to the discrepancy observed across the two studies.

Our annotators have provided valuable feedback to us, and both of them have mentioned that in some definitions, unexpected or misplaced punctuation marks or tokens occurred, which affected the overall fluency of the definition, as otherwise the definition would be considered "Very fluent". As we do not possess the original annotations, we do not know how the original authors would rate these definitions. One of the annotators also mentioned that they found the concept of fluency very ambiguous, and this may have led to confusion. From the feedback, we noticed that this annotator has also considered factuality as part of fluency, which would not happen if they were part of the original study, as we know there was an additional, separate factuality evaluation task. Yet this is an unavoidable problem since we do not know exactly what instructions have been given to the annotators, and we can only presume minimum intervention, leading to very few task instructions aside from examples. The confusion in interpreting the concept of fluency may not only lead to a lower overall score but also a lower inter-annotator agreement in the reproduction study, as the two original annotators may have reached some level of agreement on the definition of fluency, while our annotators have not.

## Acknowledgments

## 6. Bibliography

Tal August, Katharina Reinecke, and Noah A. Smith. 2022. Generating scientific definitions with controllable complexity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8298–8317, Dublin, Ireland. Association for Computational Linguistics.

Iz Beltagy, 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. Association for Computational Linguistics.

Anya Belz, Simon Mille, and David M. Howcroft. 2020. Disentangling the properties of human evaluation methods: A classification system to support comparability, meta-evaluation and reproducibility testing. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 183–194, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz, Maja Popovic, and Simon Mille. 2022. Quantified reproducibility assessment of NLP results. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16–28, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *Proceedings of*

the Fourth Workshop on Insights from Negative Results in NLP, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All that's 'human' is not gold: Evaluating human evaluation of generated text. 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 7282–7296, Online. Association for Computational Linguistics.

Lorenzo De Mattei, Michele Cafagna, Felice Dell'Orletta, and Malvina Nissim. 2020. Invisible to people but not to machines: Evaluation of style-aware HeadlineGeneration in absence of reliable human judgment. In Proceedings of the Twelfth Language Resources and Evaluation Conference, pages 6709–6717, Marseille, France. European Language Resources Association.

Lucie Flekova, Daniel Preoţiuc-Pietro, and Lyle Ungar. 2016. Exploring stylistic variation with age and income on Twitter. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 313–319, Berlin, Germany. 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. Transactions of the Association for Computational Linguistics, 9:1460–1474.

David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In Proceedings of the 13th International Conference on Natural Language Generation, pages 169–182, Dublin, Ireland. Association for Computational Linguistics.

Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. 2021. GeDi: Generative discriminator guided sequence generation. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 4929–4952, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7871–7880, Online. Association for Computational Linguistics.

Yiru Li, Huiyuan Lai, Antonio Toral, and Malvina Nissim. 2023. Same trends, different answers: Insights from a replication study of human plausibility judgments on narrative continuations. In Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems, pages 190–203, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.

Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. DExperts: Decoding-time controlled text generation with experts and anti-experts. 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 6691–6706, Online. Association for Computational Linguistics.

Ricardo Rei, 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.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval), pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Eric Smith, Orion Hsu, Rebecca Qian, Stephen Roller, Y-Lan Boureau, and Jason Weston. 2022. Human evaluation of conversations is an open problem: comparing the sensitivity of various methods for evaluating dialogue agents. In Proceedings of the 4th Workshop on NLP for Conversational AI, pages 77–97, Dublin, Ireland. Association for Computational Linguistics.

# A. HEDS Sheet

## A.1. Paper and supplementary resources

Sections 1.1–1.3 record bibliographic and related information. These are straightforward and don't warrant much in-depth explanation.

### 1.1 Details of paper reporting the evaluation experiment

#### 1.1.1 Link to paper reporting the evaluation experiment.
for preregistration. This is a reproduction experiment, and the original paper is on https://aclanthology.org/2022.acl-long.569/

#### 1.1.2 Which experiment within the paper is this form being completed for?
This form is being completed for pre-registration*
Title of experiment: Evaluating Fluency.
Section: 7 & 7.1.
Exact descriptions in Appendix A.7: "Annotators were given examples of very fluent and relevant definitions, and not at all fluent and relevant definitions before starting the task. For fluency, annotators were asked, 'How fluent is this definition?'"

### 1.2 Link to resources

#### 1.2.1 Link(s) to website(s) providing resources used in the evaluation experiment.
https://drive.google.com/drive/folders/1qq HAI_GvxO14ZoW-XGO3PMvNZnXO9mp-?usp=share_link

### 1.3 Contact details

**This part is hidden for anonymous purposes.**

#### 1.3.1 Details of the person completing this sheet

##### 1.3.1.1 Name of the person completing this sheet.
Yiru Li

##### 1.3.1.2 Affiliation of the person completing this sheet.
University of Groningen

##### 1.3.1.3 Email address of the person completing this sheet.
y.li.170@student.rug.nl

#### 1.3.2 Details of the contact author

##### 1.3.2.1 Name of the contact author.
Malvina Nissim

##### 1.3.2.2 Affiliation of the contact author.
University of Groningen

##### 1.3.2.3 Email address of the contact author.
m.nissim@rug.nl

## A.2. System Questions

Questions 2.1–2.5 record information about the system(s) (or human-authored stand-ins) whose outputs are evaluated in the Evaluation experiment that this sheet is being completed for. The input, output, and task questions in this section are closely interrelated: the value for one partially determines the others, as indicated for some combinations in Question 2.3.

### 2.1 What type of input do the evaluated system(s) take?
5. text: sentence

### 2.2 What type of output do the evaluated system(s) generate?
6. text: multiple sentences

### 2.3 How would you describe the task that the evaluated system(s) perform in mapping the inputs in Q2.1 to the outputs in Q2.2?
12. question answering

### 2.4 What are the input languages that are used by the system?
41. English

### 2.5 What are the output languages that are used by the system?
41. English

## A.3. Sample of system outputs, evaluators, and experimental design

### 3.1 Sample of system outputs
Questions 3.1.1–3.1.3 record information about the size of the sample of outputs (or human-authored stand-ins) evaluated per system, how the sample was selected, and what its statistical power is.

#### 3.1.1 How many system outputs (or other evaluation items) are evaluated per system in the evaluation experiment?
100

#### 3.1.2 How are system outputs (or other evaluation items) selected for inclusion in the evaluation experiment?
1. by an automatic random process

#### 3.1.3 Statistical power of the sample size.

##### 3.1.3.1 What method was used to determine the statistical power of the sample size?
N/A. Follow the original experiment.

### 3.1.3.2 What is the statistical power of the sample size?
N/A

### 3.1.3.3 Where can other researchers find details of the script used?
N/A

## 3.2 Evaluators
Questions 3.2.1–3.2.5 record information about the evaluators participating in the experiment.

### 3.2.1 How many evaluators are there in this experiment?
2

### 3.2.2 Evaluator Type
Questions 3.2.2.1–3.2.2.5 record information about the type of evaluators participating in the experiment.

#### 3.2.2.1 What kind of evaluators are in this experiment?
1. experts

#### 3.2.2.2 Were the participants paid or unpaid?
1. paid (monetary compensation)

#### 3.2.2.3 Were the participants previously known to the authors?
1. previously known to authors

#### 3.2.2.4 Were one or more of the authors among the participants?
2. evaluators do not include any of the authors

#### 3.2.2.5 Further details for participant type.
One participant is a non-student researcher and the other participant is a PhD student.

### 3.2.3 How are evaluators recruited?
The evaluators are recruited by in-person invitations.

### 3.2.4 What training and/or practice are evaluators given before starting on the evaluation itself?
Instructions and examples are given on the start pages of the online survey that we use to collect the results.

### 3.2.5 What other characteristics do the evaluators have?
The evaluators are expected to have high English profiency and have expertise in NLP.

## 3.3 Experimental Design
Sections 3.3.1–3.3.8 record information about the experimental design of the evaluation experiment.

### 3.3.1 Has the experimental design been pre-registered? If yes, on which registry?
2. no

### 3.3.2 How are responses collected?
Qualtrics survey.

### 3.3.3 Quality assurance
Questions 3.3.3.1 and 3.3.3.2 record information about quality assurance.

#### 3.3.3.1 What quality assurance methods are used to ensure evaluators and/or their responses are suitable?
7. None of the above. None quality assurance methods are included in the experiment, following what was in the original paper. We only made sure that the evaluators have expertise in NLP and English fluency.

#### 3.3.3.2 Please describe in detail the quality assurance methods that were used.
Expertise in NLP is expected.

### 3.3.4 Form/Interface
Questions 3.3.4.1 and 3.4.3.2 record information about the form or user interface that was shown to participants.

#### 3.3.4.1 Please include a link to online copies of the form/interface that was shown to participants.
To be determined.

#### 3.3.4.2 What do evaluators see when carrying out evaluations?
The evaluators see an information letter page which inform them of this experiment and their rights, an introduction page including examples, and then the question pages with some additional instructions.

### 3.3.5 How free are evaluators regarding when and how quickly to carry out evaluations?
3. neither of the above. We expect the evaluators to complete the whole evaluation within a set time.

### 3.3.6 Are evaluators told they can ask questions about the evaluation and/or provide feedback?
1. evaluators are told they can ask any questions during/after receiving initial training/instructions, and before the start of the evaluation

### 3.3.7 What are the experimental conditions in which evaluators carry out the evaluations?
1. evaluation carried out by evaluators at a place of their own choosing, e.g. online, using a paper form, etc.

### 3.3.8 Briefly describe the (range of different) conditions in which evaluators carry out the evaluations.
N/A

## A.4. Quality Criteria - Definition and Operationalisation

Questions in this section collect information about each quality criterion assessed in the single human evaluation experiment that this sheet is being completed for.

### 4.1 Quality Criteria

Questions 4.1.1–4.1.3 capture the aspect of quality that is assessed by a given quality criterion in terms of three orthogonal properties. They help determine whether or not the same aspect of quality is being evaluated in different evaluation experiments. The three properties characterise quality criteria in terms of (i) what type of quality is being assessed; (ii) what aspect of the system output is being assessed; and (iii) whether system outputs are assessed in their own right or with reference to some system-internal or system-external frame of reference. For full explanations see Belz et al. (2020).

#### 4.1.1 What type of quality is assessed by the quality criterion?
2. Goodness

#### 4.1.2 Which aspect of system outputs is assessed by the quality criterion?
1. Form of output

#### 4.1.3 Is each output assessed for quality in its own right, or with reference to a system-internal or external frame of reference?
1. Quality of output in its own right

### 4.2 Evaluation mode properties

Questions 4.2.1–4.2.3 record properties that are orthogonal to quality criteria (covered by questions in the preceding section), i.e. any given quality criterion can in principle be combined with any of the modes (although some combinations are more common than others).

#### 4.2.1 Does an individual assessment involve an objective or a subjective judgment
2. Subjective

#### 4.2.2 Are outputs assessed in absolute or relative terms?
1. Absolute

#### 4.2.3 Is the evaluation intrinsic or extrinsic?
1. Intrinsic

### 4.3 Response elicitation

The questions in this section concern response elicitation, by which we mean how the ratings or other measurements that represent assessments for the quality criterion in question are obtained, covering what is presented to evaluators, how they select response and via what type of tool, etc. The eleven questions (4.3.1–4.3.11) are based on the information annotated in the large scale survey of human evaluation methods in NLG by Howcroft et al. (2020).

#### 4.3.1 What do you call the quality criterion in explanations/interfaces to evaluators? Enter 'N/A' if no definition given.
Fluency

#### 4.3.2 Question 4.3.2: What definition do you give for the quality criterion in explanations/interfaces to evaluators? Enter 'N/A' if no definition given.
N/A. We provide examples though.

#### 4.3.3 Are the rating instrument response values discrete or continuous? If so, please also indicate the size.
1. Discrete
Size of the instrument: 4

#### 4.3.4 List or range of possible values of the scale or other rating instrument. Enter 'N/A', if there is no rating instrument.
1-4 Likert Scale

#### 4.3.5 How is the scale or other rating instrument presented to evaluators? If none match, select 'Other' and describe.
1. Multiple-choice options

#### 4.3.6 If there is no rating instrument, describe briefly what task the evaluators perform (e.g. ranking multiple outputs, finding information, playing a game, etc.), and what information is recorded. Enter 'N/A' if there is a rating instrument.
N/A

#### 4.3.7 What is the verbatim question, prompt or instruction given to evaluators (visible to them during each individual assessment)?
Instructions
Please read the following text and answer the questions below.
When reading definitions, please focus on their fluency. If a definition's text only says 'nan', please rate it as Not at all fluent.
Term:
Definition:
∗ How fluent is this definition?

#### 4.3.8 Form of response elicitation. If none match, select 'Other' and describe.
2. direct quality estimation

#### 4.3.9 How are raw responses from participants aggregated or otherwise processed to obtain reported scores for

**this quality criterion?**
macro-averages

**4.3.10 Method(s) used for determining effect size and significance of findings for this quality criterion.**
Pairwise independent t-tests corrected for multiple hypothesis testing using the Bonferroni-Holm correction

**4.3.11 Inter-annotator agreement**
Questions 4.3.11.1 and 4.3.11.2 record information about inter-annotator agreement.

**4.3.11.1 Has the inter-annotator agreement between evaluators for this quality criterion been measured? If yes, what method was used?**
1. yes, Krippendorff's $\alpha$

**4.3.11.2 What was the inter-annotator agreement score?**
0.45

**4.3.12 Intra-annotator agreement**
Questions 4.3.12.1 and 4.3.12.2 record information about intra-annotator agreement.

**4.3.11.1 Has the intra-annotator agreement between evaluators for this quality criterion been measured? If yes, what method was used?**
3. N/A. In our experiment, each evaluator only evaluate each item once.

**4.3.11.2 What was the intra-annotator agreement score?**
N/A

## A.5. Ethics

The questions in this section relate to ethical aspects of the evaluation. Information can be entered in the text box provided, and/or by linking to a source where complete information can be found.

**5.1 Has the evaluation experiment this sheet is being completed for, or the larger study it is part of, been approved by a research ethics committee? If yes, which research ethics committee?**
Yes. The Research Ethics Committee (CETO) of the Faculty of Arts, University of Groningen.

**5.2 Do any of the system outputs (or human-authored stand-ins) evaluated, or do any of the responses collected, in the experiment contain personal data (as defined in GDPR Art. 4, §1: https://gdpr.eu/article-4-definitions)? If yes, describe data and state how addressed.**
No. The responses are anonymized.

**5.3 Do any of the system outputs (or human-authored stand-ins) evaluated, or do any of the responses collected, in the experiment contain special category information (as defined in GDPR Art. 9, §1: https://gdpr.eu/article-9-processing-special-categories-of-personal-data-prohibited)? If yes, describe data and state how addressed.**
No

**5.4 Have any impact assessments been carried out for the evaluation experiment, and/or any data collected/evaluated in connection with it? If yes, summarise approach(es) and outcomes.**
No

# ReproHum #0124-03: Reproducing Human Evaluations of end-to-end approaches for Referring Expression Generation

**Saad Mahamood**

trivago N.V.
Düsseldorf, Germany
saad.mahamood@trivago.com

## Abstract

In this paper we describe our attempt to reproduce a single human evaluation quality criterion of the human evaluation that was in conducted in the paper "NeuralREG: An end-to-end approach to referring expression generation". In particular, this paper describes the approach and challenges involved in reproducing the human evaluation as done by the original authors of the paper, the results obtained, and what insights we have gained from attempting this particular reproduction. Insights that we hope will enable refinements to both how human evaluations are documented by author(s) and enable better reproductions of NLP experiments in the future.

**Keywords:** human evaluation, NLP, neural REG, reproduction

## 1. Introduction

There has been significant interest in understanding the issues that prevent the reproduction and repeatability of human NLP evaluations. Efforts such as the ReproHum project[1] attempts to investigate human evaluations within NLP by systematically uncovering the extent of problems of reproducibility. Uncovering these issues is especially important within the field of NLP considering the significance of human evaluations, which are seen as the "gold standard" as compared to automatic metric based evaluations, which may not correlate well with human judgement (Belz and Reiter, 2006). Past research has indicated only a minority of systems can reproduce previously reported scores and systems due either to not working non-functional code or resource limits (Belz et al., 2021b). In fact some estimates place the percentage of papers being repeatable without any significant barriers as low as 5% and at 20% if the original author(s) help is sought (Belz et al., 2023). In addition to buggy code, other issues have been observed such as flaws within the user interface to collect evaluator responses, inappropriate exclusion of evaluators and/or data points, reporting flaws, and also ethical flaws (Thomson et al., 2024).

As part of the ReproHum multi-lab study (Belz and Thomson, 2024), multiple partner labs have come have come together to reproduce existing human evaluations experiments from a chosen set of human evaluations in published NLP research papers. Papers that were vetted by the organising committee to ensure that sufficient details in terms of materials (code, data, etc.) and evalua-

tion procedures were present for a successful attempt at reproduction by a given partner lab. In addition to the original paper author(s) consent and co-operation was sought to enable the reproduction of human evaluations in their paper. Consecutively participating partner labs must follow a common reproduction approach to ensure consistency and comparability between different reproduction attempts.

This years reproduction experiment is a continuation of past years, which since 2021[2] has expanded the scope of reproduction experiments. Results from previous iterations have found the impact that different cohorts can have in the reproducibility of a given experiment (Belz et al., 2021a), or the need to lower cognitive loads for evaluators, which could potentially lead to be better reproducibility of results (Belz et al., 2022). In the 2023 edition there were three main challenges identified in trying to run reproduction results. The first was reproduction attempts encountering bugs, errors, and flaws, which were fixed differently by different reproducing authors. Secondly, reproducing authors chose different results to reproduce and report making comparability between results not possible. Finally, not all reproducing authors were able to adhere closely to the original experiment details with variations occurring such as using a different evaluation interface, or different number of evaluators (Belz and Thomson, 2023).

Based on the learnings from last year several changes have been implemented by the organis-

---

[1] ReproHum - https://reprohum.github.io

---

[2] ReproGen 2021 - https://reprogen.github.io/2021/
ReproGen 2022 - https://reprogen.github.io
ReproNLP 2023 - https://repronlp.github.io/2023

ers. There is now a revised and expanded common approach to reproduction that formalises that gives greater guidance on how the reproduction should be conducted and how the results should be reported to ensure greater comparability and standardisation between different reproduction attempts for the same paper.

In this paper we give a description of our attempt to reproduce human evaluations within the paper "NeuralREG: An end-to-end approach to referring expression generation" by Castro Ferreira et al. (2018) (section 2) and how the reproduction of the paper was conducted. We detail the challenges involved (section 3). We also detail the results obtained from the reproduction (section 4) and how they compare to the original results and the observations made by authors. Finally, we conclude with the learnings (section 5) that we have obtained based on the experiences of this reproduction experiment and describe improvements that would enable more robust reproductions of future NLP human evaluations.

## 2. Reproduction Experiment

In this reproduction experiment we were tasked with to reproduce human evaluations was "Neural-REG: An end-to-end approach to referring expression generation" (Castro Ferreira et al., 2018). The paper itself describes the creation of and an evaluation of an end-to-end neural approach for generating referring expressions, which then compared against two non-neural baseline models using the WebNLG dataset (Gardent et al., 2017). In particular, there are three neural variant systems that uses a different LSTM decoders tested by the authors and two non-neural variants:

- **OnlyNames** – A baseline non-neural model that leverages the similarity among the Wikipedia ID of an element and proper name reference to it. Basically, it replaces the underscores in a given Wikipedia ID for whitespaces.
- **Ferreira** – A second non-neural baseline model that leverages the Naive Bayes method to determine whether a given reference should be a proper name, pronoun, description, or demonstrative.
- **NeuralREG+**`Seq2Seq` – Leverages a decoding approach that models a given context vector for a given time step and concatenates the pre- and pos-context annotations averaged over time.
- **NeuralREG+**`CAtt` – A LSTM decoder that is augmented with an attention mechanism (Bahdanau et al., 2014), which for a given time step, used over the pre- and pos-context encodings.
- **NeuralREG+**`HierAtt` – Inspired by Libovický and Helcl (2017), this version implements a second attention mechanism in order to generate attention weights for pre- and pos-context summary vectors instead of concatenating them.
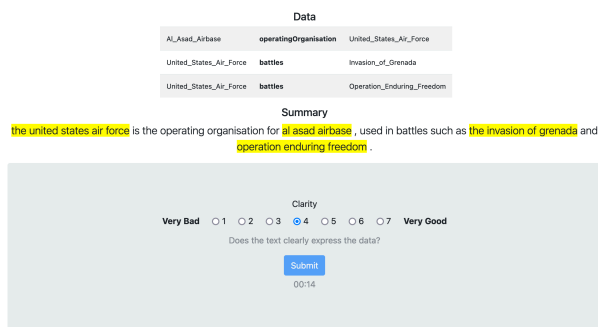


Figure 1: Evaluation interface used for rating the degree of clarity of a text containing generated referring expressions (highlighted in yellow).

Whilst these systems were evaluated using both automatic and human evaluations, the focus of this reproduction is solely on the human evaluation conducted by the original authors. In particular, the authors designed an intrinsic evaluation tasked that leveraged 24 randomly selected test WebNLG triplet instances and generated 6 target texts with referring expressions: The original (randomly selected) and five other referring expression version texts generated by each of the models described above. Using a latin square design, the authors created 144 different trials over 6 different list and designed the evaluation in a way that a given participant rated 24 trials, one for each of the 24 corpus instances, ensuring that participants saw an equal number of triplet set sizes and generated versions.

In the original experiment participants were asked to rate in a given trial three aspects for a given text containing referring expressions: *fluency*, *grammaticality*, and *clarity*. For the reproduction experiment we are tasked with only reproducing the *clarity* quality criteria aspect. Defined by the authors as whether the text clearly expresses the data. The quality criterion were rated by the participants using a seven point Likert scale. The task as done in the reproduction experiment is illustrated in figure 1, which shows a given set of triplets presented to the user in a tabular form and the text underneath with the generated referring expressions highlighted in yellow. Annotators are given 20 seconds to consider the data, the text with the generated referring expressions, and then give their ratings. This timer is unchanged from the original experiment even though the the number of quality criterion has been reduced from three to one.

Other changes to the user interface were limited in scope to accommodate ethical concerns or to update explanatory text to the fact that only one quality criteria aspect would evaluated instead of three. Changes included adding informed consent

| Aspect | Original | Reproduction |
|---|---|---|
| **Quality Criterion** | *fluency, grammaticality, clarity* | *clarity* |
| **Number of Items** | 144 | 144 |
| **Number of Systems** | 6 | 6 |
| **Number of Participants** | 60 | 60 |
| **Participants per Item** | 10 | 10 |
| **Items per Participant** | 24 | 24 |
| **Recruitment Platform** | *Amazon MTurk* | *Prolific* |
| **Compensation** | *unknown* | *£12.00 per hour equivalent* |
| **Participation controls** | *unknown* | *none* |
| **Age** | *Average 36 years* | *Majority 18-24 years (43%)* |
| **Gender Split** | *27 females, 33 males* | *35 females, 25 males* |
| **English Proficiency** | *Native: 44, Fluent: 14, Basic: 2* | *Native: 37, Fluent: 21, Basic: 2* |

Table 1: Methodological similarities & differences between the original and reproduction human evaluations.

## 3. Methodology & Challenges

Participants for the original experiment were recruited from Amazon Mechanical Turk, with 60 participants recruited and 10 assigned for each of the six lists. In the reproduction experiment, participants were recruited instead from Prolific[3] in agreement with the ReproHum organisers to ensure every reproduction experiment used a standardised crowd working platform. Whilst, the original experiment does not detail the degree of compensation given to participants, for the reproduction experiment participants were paid the equivalent of the UK living wage[4] of £12.00 per hour for their participation. Table 1 details the methodological and participatory similarities and differences between the two experiments. In terms of demographics, in the reproduction experiment the age is much younger that in the original experiment with 43% of participants reporting themselves to be between 18-24 years old and there is a greater proportion of participants identifying as female compared to the original experiment. For English proficiency, the distribution between the original and reproduction are fairly similar although with a slight more number of fluent instead of native English speakers.

The experimental data and user interface was taken from the original published source code repository[5]. The main evaluation interface, was a web application that was written in PHP with the purpose of handling collecting user responses and assigning users to an equal number of evaluation lists. However several challenges were encountered in attempting to reuse the original experimental data and user interface:

- The database structure was not available in the GitHub repository. As part of the reproduction this had to be recreated by interpreting the existing PHP code and through trial and error.
- The order of evaluations items was not defined for each list as this was was encoded in the not provided database.
- Lack of detailed version information for both the software used for the evaluation interface and the analysis code.

For the second point, whilst the start item for each of six lists was hard coded into the PHP code the order of subsequent items was not known. Therefore in coordination with the ReproHum organisers it was decided to randomise the order of items for each of the six lists. However, this change may have lead to a potential deviation from how the original experiment was conducted by the authors. Whilst, writing this report it was discovered that the code for the generate the trial lists was hidden in a python file that was used for computing the result statistics of the human evaluations.

In addition to setting up the reproduction experiment by using the original experiment's codebase a Human Evaluation Datasheet (HEDS) (Shimorina and Belz, 2022) was also completed[6]. The HEDS form records in a standardised way the properties of human evaluations to support comparability, meta-evaluation, and reproducibility of human evaluations.

## 4. Results

In the original experiment the authors made the following observations with respect to how the neu-

---

[3]Prolific - https://www.prolific.com

[4]UK Living Wage - https://www.livingwage.org.uk

[5]NeuralREG - https://github.com/ThiagoCF05/NeuralREG

[6]ReproNLP 2024 HEDS forms - https://github.com/nlp-heds/repronlp2024

| | Original | Reproduction | CV* |
|---|---|---|---|
| *OnlyNames* | 4.90 | 4.92 | 0.4061121348816013 |
| *Ferreira* | 4.93 | 4.69 | 4.974662575306527 |
| NeuralREG+Seq2Seq | 4.97 | 4.97 | 0.0 |
| NeuralREG+CAtt | 5.26 | 4.97 | 5.652620418943544 |
| NeuralREG+HierAtt | 5.13 | 5.04 | 1.7646111347510636 |
| *Original* | 5.42 | 5.22 | 3.7481401922344113 |

Table 2: Clarity mean average results from both original and reproduction human evaluation. Unbiased coefficient of variation values (**CV\***) calculated using the definition by Belz (2022). Original results are from (Castro Ferreira et al., 2018).

ral models performed against the baseline models and the original texts:

1. *"...all three neural models scored higher than the baselines on all metrics, with especially Neural-REG+CAtt approaching the ratings for the original sentences."*

2. *"...differences between the neural models were small"*

3. *"The results for the 3 different decoding methods of NeuralREG also did not reveal a significant difference."*

4. *"...the original texts were rated significantly higher than both baselines in terms of the three metrics...and than NeuralREG+Seq2Seq in terms of clarity."*

From the results of the reproduction the claims made by the original authors do all hold up and are backed by the results as shown in table 2. This table also includes a column for coefficient of variation for small sample sizes using the methodological approach defined by Belz (2022). Correlations between the original and the reproduction results using both Pearson's *r* of *r=0.783* (*p=0.065*) and Spearman's $\rho$ of $\rho$=*0.840* (*p=0.036*) were calculated, with both showing statistically significant positive correlations.

With the exception of the *OnlyNames* (slight improvement over original) and the Neural-REG+Seq2Seq systems (same result as original) all other variants showed a decrease in average clarity ratings as compared to the original evaluation. One interesting result is that of the NeuralREG+CAtt system, which showed a marked decrease. Nevertheless, the system still performed as equally as well as the Neu-ralREG+Seq2Seq and better than the baseline non-neural systems. One possible explanation for the the observed differences could be due to the different cohort of evaluators in the reproduction as compared the original study. The evaluators in the reproduction are much younger and have a greater degree of English language proficiency and this may have lead to the observed variances seen.

## 5. Conclusion

In this paper we have conducted a successful reproduction of the results obtained in the original human evaluation by Castro Ferreira et al. (2018). There was slight variances in the reported scores in the reproduction, which for a majority of them had slightly lower scores that those originally reported with the exception of two of variants. However, the finding by that the original authors that the neural systems outperform the baselines, whilst underperforming the original text variant holds true and is reconfirmed in this reproduction. In addition, the results obtained in the reproduction show statistically significant positive correlations against the original results.

There are several factors that may have led to this reproduction to having a successful outcome. Factors such as the completeness of the web interface code, the presences of both original collected dataset, and the presence of functional analysis code. Nevertheless, there are areas of improvements. Such as having complete documentation for setting up the experiment. For example, the issue with respect to the order items for each of the six lists could have been mitigated with documentation by the original study authors on the places to look when trying to recreate a given study. Additionally, better documentation would help to remove uncertainty in two aspects. The controls applied for recruiting participants (if any) and the versions of software and libraries used for both the web interface and analysis code. Finally, the missing database schema could of hindered the reproduction experiment from being run at all, but thankfully was worked around with some reverse engineering of the web interface code. Incorporate these improvements would would not only reduce uncertainty, but also reduce the friction in trying to attempt a reproduction by future prospective reproducing authors.

## Acknowledgements

# 6. Bibliographical References

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv preprint arXiv:1409.0473*.

Anja Belz and Ehud Reiter. 2006. Comparing automatic and human evaluation of NLG systems. In *11th conference of the European chapter of the association for computational linguistics*, pages 313–320.

Anja Belz, Anastasia Shimorina, Shubham Agarwal, and Ehud Reiter. 2021a. The ReproGen shared task on reproducibility of human evaluations in NLG: Overview and results. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 249–258.

Anja Belz and Craig Thomson. 2023. The 2023 ReproNLP Shared Task on Reproducibility of Evaluations in NLP: Overview and Results. In *Proceedings of the 3rd Workshop on Human Evaluation of NLP Systems*, pages 35–48.

Anya Belz. 2022. A Metrological Perspective on Reproducibility in NLP*. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Shubham Agarwal, Anastasia Shimorina, and Ehud Reiter. 2021b. A systematic review of reproducibility research in natural language processing. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 381–393, Online. Association for Computational Linguistics.

Anya Belz, Anastasia Shimorina, Maja Popović, and Ehud Reiter. 2022. The 2022 ReproGen shared task on reproducibility of evaluations in NLG: Overview and results. In *Proceedings of the 15th International Conference on Natural Language Generation: Generation Challenges*, pages 43–51, Waterville, Maine, USA and virtual meeting. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2024. The 2024 ReproNLP Shared Task on Reproducibility of Evaluations in NLP: Overview and Results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, and Simon Mille. 2023. Non-repeatable experiments and non-reproducible results: The reproducibility crisis in human evaluation in NLP. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3676–3687, Toronto, Canada. Association for Computational Linguistics.

Thiago Castro Ferreira, Diego Moussallem, Ákos Kádár, Sander Wubben, and Emiel Krahmer. 2018. NeuralREG: An end-to-end approach to referring expression generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1959–1969, Melbourne, Australia. Association for Computational Linguistics.

Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. Creating training corpora for NLG microplanners. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 179–188, Vancouver, Canada. Association for Computational Linguistics.

Jindřich Libovický and Jindřich Helcl. 2017. Attention strategies for multi-source sequence-to-sequence learning. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 196–202, Vancouver, Canada. Association for Computational Linguistics.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Craig Thomson, Ehud Reiter, and Anya Belz. 2024. Common Flaws in Running Human Evaluation Experiments in NLP. *Computational Linguistics*, pages 1–11.

# ReproHum #0087-01: Human Evaluation Reproduction Report for *Generating Fact Checking Explanations*

**Tyler Loakman[1], Chenghua Lin[2]**
[1]Department of Computer Science, The University of Sheffield, UK
[2]Department of Computer Science, The University of Manchester, UK
tcloakman1@sheffield.ac.uk, chenghua.lin@manchester.ac.uk

## Abstract

This paper presents a partial reproduction of *Generating Fact Checking Explanations* by Atanasova et al. (2020) as part of the ReproHum (Belz and Thomson, 2024) element of the ReproNLP shared task to reproduce the findings of NLP research regarding human evaluation. This shared task aims to investigate the extent to which NLP as a field is becoming more or less reproducible over time. Following the instructions provided by the task organisers and the original authors, we collect relative rankings of 3 fact-checking explanations (comprising a gold standard and the outputs of 2 models) for 40 inputs on the criteria of *Coverage*. The results of our reproduction and reanalysis of the original work's raw results lend support to the original findings, with similar patterns seen between the original work and our reproduction. Whilst we observe slight variation from the original results, our findings support the main conclusions drawn by the original authors pertaining to the efficacy of their proposed models.

## 1. Introduction

Recently, many works have investigated the role of human evaluation in assessing the quality of outputs in the field of Natural Language Processing (NLP) and Natural Language Generation (NLG) (Belz et al., 2023; Clark et al., 2021; van der Lee et al., 2019). Whilst human evaluation is often seen as the gold standard method of evaluation which takes into account the perceptions of real human end-users, there is much debate over the reproducibility of such evaluation (Belz et al., 2023; Howcroft et al., 2020). Automatic metrics, whilst scalable, frequently demonstrate poor concurrent validity, correlating poorly with human judgements (Liu et al., 2024; Zhao et al., 2023; Alva-Manchego et al., 2021; Reiter, 2018; Belz and Reiter, 2006). However, the performance of human evaluation has likewise been shown to have multiple flaws, including ill-defined evaluation criteria compounded by the absence of sufficient evaluator/annotator training to attenuate the subjectivity of the texts being rated from the subjective interpretation of the evaluation criteria itself. Furthermore, several works have discussed the presence of poorly selected human panels, including sufficient language proficiency and task understanding Schoch et al. (2020). This is further hindered by the choice of many works to obfuscate these shortcomings by neglecting to report any demographic information regarding participants, including for highly subjective language types such as humour (Loakman et al., 2023). Such discrepancies have resulted in widespread troubles in reproducing the results of different works in NLP (Thomson et al., 2024).

It is for reasons such as these that the ReproHum shared task aims to shine a spotlight on the level of reproducibility within the field of NLP through the mass reproduction of contemporary research through its many partner labs so that poor practices are identified and a record can be made of the progress of reproducibility over time, as researchers become increasingly aware of the best practices to follow in performing human evaluation in their works.

## 2. Background

As participants in the ReproHum project, we selected the paper *Generating Fact Checking Explanations* by Atanasova et al. (2020) as the focus of our reproduction, owing to interest in the topic of explanation generation, and previous experience of being part of evaluator panels for similar research. Through the automatic selection process, the ReproHum team identified the single experiment and criterion that we were to attempt to reproduce the results from, as introduced in §4.

Owing to our participation in the ReproHum project (Belz and Thomson, 2024), we were provided with the following materials: (i) a guide to the common approach to reproduction, (ii) the original paper and dataset required to perform a reproduction, and (iii) additional documents pertaining to clarifications and additional information provided by the original authors once contacted. During this process, the authors of this paper (and therefore the team performing the reproduction) did not con-

tact the authors of the original work directly at any stage.

In performing this reproduction, we adhered to the following criteria outlined in the documentation provided by the ReproHum organisers. All participants were paid minimally to the UK National Living Wage (12GBP per hour) as set by the Repro-Hum team for pair pay, in which we specifically paid 15GBP for this task and paid via Amazon Vouchers from our estimation that the task would take approximately 1.25hrs (which was confirmed by our evaluators following completion). Additionally, this work underwent ethical review and approval by the ethics review board of the primary author's institution (where all participants in this reproduction were also selected).

## 3. Original Study

In recent years with the widespread sharing of misinformation and the coining of "fake news", the need for accurate and reliable fact-checking systems has grown exponentially. While existing systems have demonstrated impressive performance, their "black box" nature often obscures the reasoning behind their predictions. This lack of transparency can hinder user trust and limit the adoption of these systems. Atanasova et al. (2020) identified an overall research focus on the veracity prediction task of news claims in existing research and a lack of work focusing on generating natural language explanations to justify these veracity predictions. They aimed to address the main drawback of a black-box system by generating explanations to support the assigned veracity labels. To do this, the authors leverage detailed fact-checking reports (termed "ruling comments") published alongside veracity labels by fact-checking organisations to produce explanations that resemble human-written justifications. This approach is further bolstered through a multi-task learning framework, where explanation generation is jointly optimised with a veracity prediction task for a DistilBERT (Sanh et al., 2020) based model. This joint training enables the system to identify regions in the ruling comments that are not only close to the gold standard explanation but also contribute to the overall fact-checking decision.

### 3.1. Evaluation

The authors evaluate their approach using both automatic and human evaluation methods. While automatic evaluation relies on the standard metric of ROUGE (Lin, 2004), human evaluation focuses on a range of different criteria listed below, alongside their original definitions:

- **Coverage** - The explanation contains important, salient information and does not miss any important points that contribute to the fact check.

- **Non-redundancy** - The summary does not contain any information that is redundant/repeated/not relevant to the claim and the fact check.

- **Non-contradiction** - The summary does not contain any pieces of information that are contradictory to the claim and the fact check.

Based on these criteria, evaluators are requested to rank different explanations based on their performance on each criterion (as well as providing an *Overall* ranking). The original results in Atanasova et al. (2020) demonstrate that the multi-task learning approach leads to improved performance for both veracity prediction and explanation generation. Notably, the generated explanations achieve better coverage and overall quality compared to explanations trained solely to mimic human justifications. This suggests that the joint training framework allows the system to capture the knowledge required for accurate fact-checking, leading to more informative and relevant explanations. In our reproduction, we focus solely on the underlined criterion of *Coverage*.

## 4. Reproduction Setting

**Task Setting** As directed by the ReproHum team, we performed our reproduction on a single element of the original work by Atanasova et al. (2020) regarding evaluating outputs on the aforementioned criteria of *Coverage*. We presented the same instructions to participants as presented by Atanasova et al. (2020) with minor changes, as presented in Figure 1. These changes exclusively involve the removal of information regarding other evaluation criteria used in the original study outside of *Coverage*, including *Non-redundancy*, *Non-contradiction*, and a holistic *Overall* rating. We additionally remove all mention of the separate Task 2 which is not the subject of this reproduction. As with the original study, we performed our reproduction experiment by having participants place their relative preference rankings of 3 systems (i.e., a gold standard and two models) in a spreadsheet facilitated via Google Sheets. Within this, 3 columns follow the 3 explanations (from the 3 different models) to place rankings (where the $n$-th column contains the ranking for the $n$-th justification), as outlined in Figure 1. In line with the recommended approach to performing reproductions presented by the ReproHum team, we additionally incorporate data validation techniques in the form of drop-down

Figure 1: Modified instructions from Atanasova et al. (2020) presented to participants within the reproduction. we made minor modifications to the original instructions presented to participants in order to remove information related to tasks and criteria that were not to be assessed in this reproduction.

boxes containing rankings of 1-3. This ensured that participants only entered valid options in the ranking task. We present model outputs to participants in the same shuffled order presented in the original paper to also avoid order effects and bias towards particular columns. In total, each participant annotated 120 items, consisting of the outputs of 3 systems (including the human gold standard) for 40 inputs. We also make available a HEDS datasheet (Shimorina and Belz, 2022) detailing the process of our reproduction study.[1]

**Evaluator Demographics**   In the original work by Atanasova et al. (2020) we have limited demographic details regarding the participants. However, we are aware that they are colleagues of the authors and have experience in fact-checking annotation tasks, whilst not exclusively being native speakers of the target language. In our replication, we use 3 Ph.D. students in Natural Language Processing, all of which have experience in fact-checking and

related tasks (e.g., misinformation/rumour detection). All participants in our reproduction also have a professional working level of English fluency.

## 5.  Results

We present the results of the original study and our reproduction in Table 1. Due to minor discrepancies in the specific evaluated materials (owing to some evaluators in the original work assessing approximately 80 items, and others assessing only 39, with some omissions), we additionally report what we term a "recreation", where we reanalyse the original paper's raw data to facilitate a direct comparison against only the same 40 inputs as presented to our evaluators. In the original work by Atanasova et al. (2020), the criterion of *Coverage* is shown to have low inter-annotator agreement as calculated via Krippendorff's Alpha (Krippendorff, 2019), reporting $\alpha = 0.26$ across their 3 evaluators. In our reproduction, we find slightly better agreement among our participants, with $\alpha = 0.35$ when specifically accounting for an ordinal level of measurement, whilst we find agreement across the 40

---

[1]Available at https://github.com/nlp-heds/repronlp2024.

| Original | | | |
|---|---|---|---|
| **Annotators** | **Gold** | **Explain-Extr** | **Explain-MT** |
| All | **1.48** | 1.89 | <u>1.68</u> |
| 1st | **1.50** | 2.08 | <u>1.87</u> |
| 2nd | **1.74** | 2.16 | <u>1.84</u> |
| 3rd | **1.21** | 1.42 | <u>1.34</u> |
| CV* | 9.00% | 8.10% | 5.76% |
| **Recreation** | | | |
| **Annotators** | **Gold** | **Explain-Extr** | **Explain-MT** |
| All | **1.52** | 1.87 | <u>1.66</u> |
| 1st | **1.55** | 2.05 | <u>1.85</u> |
| 2nd | 1.82 | 2.15 | <u>**1.77**</u> |
| 3rd | **1.18** | 1.41 | <u>1.36</u> |
| CV* | 6.35% | 9.16% | 6.96% |
| **Reproduction** | | | |
| **Annotators** | **Gold** | **Explain-Extr** | **Explain-MT** |
| All | **1.62** | 2.05 | <u>1.78</u> |
| 1st | **1.60** | 2.30 | <u>2.03</u> |
| 2nd | 1.60 | 1.86 | <u>**1.55**</u> |
| 3rd | **1.65** | 1.98 | <u>1.75</u> |

Table 1: Comparison between Atanasova et al. (2020) and our reproduction on the criterion of "Coverage". Values present the Mean Average Ranks (MAR) of the explanations. The explanations come from the gold justification (**Gold**), the generated explanation (**Explain-Extr**), and the explanation learned jointly (**Explain-MT**) with the veracity prediction model. A lower MAR indicates a better average ranking. For each row, the best results are in **bold**, and the best automatically generated explanations are <u>underlined</u>. "Annotators" refers to each individual rater, whilst "All" is the mean across all annotators. *CV\** refers to the Coefficient of Variation for the mean ratings of the 3 systems compared to our reproduction results following the implementation by Belz (2022). *Original* refers to the results presented in the original paper by Atanasova et al. (2020), whilst *Recreation* refers to the results we gain by reanalysing the original study's data exclusively for the same sample that our evaluators were presented. Finally, *Reproduction* refers to the results of our reproduction study using our new evaluators. The ordering of annotators across *Recreation* and *Original* should be considered arbitrary, as we cannot guarantee each line corresponds to the same annotator as the original.

evaluated inputs in the original data to be very similar to what was reported for the particular subset used by the authors in the original work ($\alpha = .27$)

In terms of overall patterns seen in the data, the results of our reproduction can be seen to differ slightly from those of the original in terms of overall rankings. Firstly, in the original study, the golden human-authored explanations were preferred by all participants, whilst this is not seen to be the case in our reproduction or in our reanalysis of a specific subset of the original paper's raw data (i.e., *recreation*) Instead, we find only 2 of our 3 participants to rank the golden explanations in their expected 1st place. However, in terms of the automatically generated explanations we observe the *Explain-MT* model (where the explanation is learnt jointly with the veracity prediction model) to outperform *Explain-Extr* (where the auxiliary veracity prediction model is learnt separately), mirroring the results presented in the original work.

Furthermore, when aggregating the results of all 3 evaluators in our reproduction, we can see that the overall rankings assigned to each output are higher (i.e., worse) than the findings of Atanasova et al. (2020). However, whilst our raw figures differ from the original findings (owing to the relatively subjective task criteria and small evaluator panel sizes), our findings reflect the same overall patterns as the original work, with the human-authored golden explanations *Gold* outperforming the authors' proposed models in the majority of cases, whilst the more complex *Explain-MT* model, which is trained alongside a veracity prediction task, outperforms the *Explain-Extr* model that learns to generate explanations in isolation.

To compare against the original study's findings, we calculate correlations between our results and those provided by the original paper's authors using Spearman's $\rho$ and Pearson's $r$. Due to the original work's raw data having results for more than 40

trials, and with some missing values, we assess only the same 40 trials as presented to our participants (equivalent to the *Recreation* in Table 1) and calculate the mean rank given to each output by the evaluators (which is robust to cases where not all evaluators in the original work assessed a given output). The results show a strong correlation between the results of our reproduction and the original study ($\rho = .524$ and $r = .541$, which are both significant at $\alpha = .01$), demonstrating that we were able to reproduce the general evaluator preferences observed in the original experiment.

## 6. Conclusion

In this paper, we have presented our reproduction findings for an element of human evaluation presented in Atanasova et al. (2020) regarding the criteria of *Coverage* to compare gold standard fact-checking explanations with 2 proposed models. In terms of overall comparison with the original work, we find a higher level of rating agreement among our evaluator panel than demonstrated in the original work but also observe a slightly different overall pattern than presented by the original authors, with one of the proposed models ranking higher than the gold standard human-authored explanation from 1 of our 3 participants. We do, however, observe the same pattern when reanalysing the raw data from the original study, focussing exclusively on the same subset of examples presented to our evaluators in the reproduction. Additionally, our reproduction lends credence to the results presented by Atanasova et al. (2020) regarding the model trained to generate explanations alongside a veracity prediction model (Explain-MT) outperforming the model that is trained to generate explanations in isolation (Explain-Extr) in terms of human rankings. It is important to note, however, that the result of our reproduction covers only one of the multiple human evaluation criteria on which the raters were asked to assess the generations in the original work, and this pattern may not necessarily be present across all different criteria.

Overall, we reiterate the importance of performing reproduction studies such as this in order to assess the trend of reproducibility within the field of NLP. Within this paper, we have successfully reproduced the findings of the original work with some minor variability (likely owing to the small size of the evaluation panels in the original work, and consequently our reproduction). This is particularly salient for the topic of generating fact-checking explanations that Atanasova et al. (2020) tackle, as this constitutes a high-impact application of NLP with an increased need for reliable and robust models and evaluation procedures in order to avoid the effects of misinformation.

## 7. Acknowledgements

## 8. Bibliographical References

Fernando Alva-Manchego, Carolina Scarton, and Lucia Specia. 2021. The (un)suitability of automatic evaluation metrics for text simplification. *Computational Linguistics*, 47(4):861–889.

Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020. Generating fact checking explanations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7352–7364, Online. Association for Computational Linguistics.

Anja Belz and Ehud Reiter. 2006. Comparing automatic and human evaluation of NLG systems. In *11th Conference of the European Chapter of the Association for Computational Linguistics*, pages 313–320, Trento, Italy. Association for Computational Linguistics.

Anya Belz. 2022. A metrological perspective on reproducibility in NLP*. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der

Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *Proceedings of the Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All that's 'human' is not gold: Evaluating human evaluation of generated text. 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 7282–7296, Online. Association for Computational Linguistics.

David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 169–182, Dublin, Ireland. Association for Computational Linguistics.

Klaus Krippendorff. 2019. Content analysis: An introduction to its methodology.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Yiqi Liu, Nafise Sadat Moosavi, and Chenghua Lin. 2024. Llms as narcissistic evaluators: When ego inflates evaluation scores.

Tyler Loakman, Aaron Maladry, and Chenghua Lin. 2023. The iron(ic) melting pot: Reviewing human evaluation in humour, irony and sarcasm generation. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6676–6689, Singapore. Association for Computational Linguistics.

Ehud Reiter. 2018. A structured review of the validity of BLEU. *Computational Linguistics*, 44(3):393–401.

Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter.

Stephanie Schoch, Diyi Yang, and Yangfeng Ji. 2020. "this is a problem, don't you agree?" framing and bias in human evaluation for natural language generation. In *Proceedings of the 1st Workshop on Evaluating NLG Evaluation*, pages 10–16, Online (Dublin, Ireland). Association for Computational Linguistics.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Craig Thomson, Ehud Reiter, and Anya Belz. 2024. Common Flaws in Running Human Evaluation Experiments in NLP. *Computational Linguistics*, pages 1–11.

Chris van der Lee, Albert Gatt, Emiel van Miltenburg, Sander Wubben, and Emiel Krahmer. 2019. Best practices for the human evaluation of automatically generated text. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 355–368, Tokyo, Japan. Association for Computational Linguistics.

Kun Zhao, Bohao Yang, Chenghua Lin, Wenge Rong, Aline Villavicencio, and Xiaohui Cui. 2023. Evaluating open-domain dialogues in latent space with next sentence prediction and mutual information. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 562–574.

# ReproHum #0892-01: The painful route to consistent results: A reproduction study of human evaluation in NLG

**Irene Mondella⋆, Huiyuan Lai◇, Malvina Nissim◇**

⋆ILC-CNR / University of Pisa, Italy
◇CLCG, University of Groningen, the Netherlands
i.mondella@studenti.unipi.it
{h.lai,m.nissim}@rug.nl

## Abstract

In spite of the core role human judgement plays in evaluating the performance of NLP systems, the way human assessments are elicited in NLP experiments, and to some extent the nature of human judgement itself, pose challenges to the reliability and validity of human evaluation. In the context of the larger ReproHum project, aimed at running large scale multi-lab reproductions of human judgement, we replicated the understandability assessment by humans on several generated outputs of simplified text described in the paper "Neural Text Simplification of Clinical Letters with a Domain Specific Phrase Table" by Shardlow and Nawaz, appeared in the Proceedings of ACL 2019. Although we had to implement a series of modifications compared to the original study, which were necessary to run our human evaluation on exactly the same data, we managed to collect assessments and compare results with the original study. We obtained results consistent with those of the reference study, confirming their findings. The paper is complete with as much information as possible to foster and facilitate future reproduction.

**Keywords:** human evaluation, reproducibility, ReproHum

## 1. Introduction

Human evaluation of model performance plays a central role in Natural Language Processing (NLP). This is particularly true in the broadly defined area of Natural Language Generation (NLG), which encompasses machine translation, rephrasing, summarisation, etc, i.e., any modelling task whose output consists in some generated text. Indeed, the large variability in acceptable outputs does not allow for an exhaustive set of gold references to be pre-produced, as is instead the case for classification tasks. For the same reason, automatic metrics must be used that are able to capture some degree of similarity between references and different but potentially valid outputs, and cannot exploit an exact correspondence of reference and output.

Developments in NLG evaluation have seen the direct incorporation of human judgements into trainable metrics, such as COMET (Rei et al., 2020), leading to much higher correlations to human assessments. While on the one hand the development of metrics that better align to human judgement appears to be a very promising direction, on the other hand the optimism could be tainted by findings along another avenue of research, dedicated to the *reproducibility* (and therefore reliability) of human judgement.

Recent efforts conducted in the context of the ReproGen shared evaluation campaigns (Belz et al., 2021, 2022) and especially the preliminary findings of ReproHum[1] (Belz et al., 2023), a cooperative project aimed to test the replicability of human eval-

uations reported in existing NLP papers through large-scale reproductions across multiple research groups, have shed some worrying light on the reliability – and thus validity – of human assessments themselves. Strikingly, Belz et al. (2023, p. 5) report "that only a small fraction of previous human evaluations in NLP can be repeated under the same conditions, hence that their reproducibility cannot be tested by repeating them."

The present paper reports on a reproduction experiment which is also part of the ReproHum project (Belz and Thomson, 2024), as an ongoing effort to further explore the extent to which human judgements elicited in NLP, and in this context more specifically NLG experiments, can be considered reliable and what mostly affects reproduction. As part of ReproHum, our work follows the research template provided by the project coordination team; this paper presents our results accordingly, thus following specific guidelines and reporting templates. The experiment was pre-registered through the Human Evaluation Data Sheet (HEDS[2]) as proposed by Shimorina and Belz (2022). We first introduce the details of the original experiment and the human evaluation included therein, and then describe our own reproduction study, specifically focusing on all the adjustments we had to made in our experiments compared to the original evaluation setup. We compare results critically, running a correlation analysis and comparing inter-annotator agreement across the two studies. We observe that our efforts in faithfully reproducing the original human evalua-

---

[1] https://reprohum.github.io/

[2] Details at the following link: https://github.com/nlp-heds/repronlp2024

| Complex Term | Simple Term |
|---|---|
| ability to be ambulant | ability to walk |
| carcinoma of stomach | cancer of stomach |
| hypertension | high blood pressure |
| osteophyte | bony spur |

Table 1: Examples of term pairs for phrase table.

tion have, in this case, brought promising results: the findings of the reference paper were confirmed by our reproduction, and the changes we had to make to the original experimental design did not affect the consistency between the two studies.

## 2. Overview of Original Study

We aim to reproduce the human evaluation experiment of text simplification in "Neural Text Simplification of Clinical Letters with a Domain Specific Phrase Table" by Shardlow and Nawaz (2019). Text simplification is the process of automatically paraphrasing a text to improve its understandability while preserving its original meaning (Al-Thanyyan and Azmi, 2021). This has a wide range of applications, such as helping non-native speakers and bridging the gap between layman and expert.

### 2.1. Task and Models

This original study aims to use text simplification methods to automatically aid patient understanding of clinical letters containing complex medical terminology (see examples in Table 5). Specifically, based on the SNOMED-CT clinical thesaurus (Donnelly, 2006), the authors created a phrase table that links complex medical terminology to simpler vocabulary (see Table 1), which is used to augment existing neural text simplification systems. To assess the impact of the proposed method on the ease of understanding sentences, human judgment is elicited to evaluate three different systems as well as the original sentences, for a total of four versions of the same sentence:

- **Original Texts (ORIG):** The original texts appear after preprocessing, which ensures that they are equivalent to the transformed texts and that any effects would be from the simplification system, not the preprocessing.

- **NTS:** The original sentences were modified by the Neural Text Simplification (NTS) system (Nisioi et al., 2017), which uses the open-source OpenNMT (Klein et al., 2017) library that provides sequence to sequence learning between a source and target language.

- **NTS + Phrase Table (NTS + PT):** The original sentences were modified by NTS, but when

OpenNMT identified a word as being out-of-vocabulary, this system (the one proposed by the authors of the original paper) will use the phrase table to replace it.

- **Phrase Table Baseline (PTB):** To demonstrate the advantages of using the phrase table in tandem with the NTS system, the proposed baseline is to only apply the phrase table to every word that could be replaced in the text.

The simplified sentences, generated by the systems described above, as well as the original version, are assessed by means of human evaluation.

### 2.2. Human Evaluation

The original study selected 50 source texts from two different datasets: i2b2 (Uzuner et al., 2007), which is a dataset of 899 discharge summaries, and MIMIC-III v1.4 (Johnson et al., 2016), which contains over 58,000 hospital records, with detailed clinical information. In this way, they obtained 100 instances: for each of them, 3 different simplified versions were created using the methods described in subsection 2.1, obtaining 100 4-tuples of parallel sentences. Texts within a 4-tuple are identical except for the modifications made by each system. No two sentences in a 4-tuple are the same.

The human evaluation was conducted on Figure Eight, a crowd-sourcing platform that no longer exist. Each 4-tuple has been assessed by 10 annotators, and each annotator could complete a maximum of 20 annotations, with the aim of obtaining a wide variety of perspectives on the data. No annotator saw the same 4-tuple twice.

To ensure the quality of annotations, workers with a higher than average rating on the Figure Eight platform were selected (level 2 and above), and a set of test annotations was designed to filter out bad-actors. From the analysis of the raw results, we found that there was a total of 8 test annotations, and most of the participants had to answer to all of them.

For each 4-tuples, annotators have been asked to rank the 4 sentences according to their ease of understanding, where the top-ranked sentence (rank 1) is the easiest to understand, while the bottom-ranked sentence (rank 4) is the hardest. Furthermore, it was specified that, in the case of 2 sentences of equal complexity, the annotator should order them according to the order of presentation. In total, 1000 annotations (100 instances with 10 annotations each) were collected. However, 20 of them were identified as not using all 4 ranks, i.e. 2 or more sentences were at the same ranking level. In these cases, the specific annotation was removed in the final analysis, resulting in 980 rankings.

| Setting | Original Study | Replicated Study |
|---------|---------------|------------------|
| Platform | Figure Eight | Prolific |
| Participants | 98 | 40 |
| Conditions | $\geq$level 2 | acceptance rate$\geq$99% & completed tasks $\geq$ 200 region filter: UK, USA, Australia, Canada |
| Filtering | a set of test annotations | 3 additional test annotations |
| Reward | Unknown | £12 per hour |

Table 2: Human evaluation settings in original and replicated study.

Finally, the authors design a metric to calculate the average rank $r_s$ of a system $s$, which is described in Equation 1.

$$r_s = \frac{\sum_{i=1}^{4} i \times f(s, i)}{\sum_{i=1}^{4} f(s, i)} \quad (1)$$

where $i$ is a rank from 1 to 4 and $f(s, i)$ is a function that maps the system and rank to the number of times that system is placed at that rank.

## 3. Reproduction Study

In our reproduction study, we strictly followed the settings of the human evaluation performed by the authors of the original work, although some adjustments had to be made for various reasons.

First, we couldn't use the crowd-sourcing platform used in the original study, because it doesn't exist anymore, so we used instead Prolific[3]. One of the main differences between these two platforms is that in Prolific it is necessary to set in advance the number of items to be evaluated by each participant. Analysing the raw results of the original paper, we assume that this constraint was not present in Figure Eight, since 76 participants evaluated 20 4-tuples (the maximum number set by the authors of the original study, included the test annotations), and 22 participants rated fewer items. In total, in the original evaluation, 98 annotators were recruited. In our case, however, it was necessary to create surveys of a fixed length. To conform our reproduction to the experimental design adopted in the ReproHum project, we created surveys containing 25 instances. Also, to ensure quality of annotations, we added 3 additional test annotations to each survey to filter out bad actors. Since the total number of instances is 100, this resulted in 4 different surveys, each of them presented to 10 different participants, for a total of 40 annotators. We made sure that no annotator participated in more than one survey.

Another difference in our replication, made necessary by the use of a different crowd-sourcing platform, regards the selection criteria for participants. Since we do not know how the participants' rating was calculated in Figure Eight, we opted to set, on Prolific, a minimum acceptance rate of 99% and a minimum completed tasks of 200. In addition, we saw from the original raw results that all the evaluators were in the United Kingdom, United States, or Australia. Whether this is by design or not we cannot tell for sure; however, because of this strong evidence, we set a region filter on English-speaking countries UK, US, Australia and Canada.

Another point on which we acted independently of the original experiment is the compensation due to the annotators, not specified in the reference paper. In our reproduction, we followed the current UK minimum wage of £12 per hour, following the general recommendation of the ReproHum project. Estimating a minimum completion time of 30 minutes per survey, we payed £6 per participant.

The differences in settings between the human evaluation performed by us and the original one are summarized in Table 2.

A screenshot of the annotation interface we created is shown in Figure 1, with instructions reported also as a screenshot in Figure 2 (the latter in the Appendix). Instead of creating a different question for each of the 4 sentences, as in the original annotation interface, we opted for a drag-and-drop system, that we find more intuitive. The instruction page, on the other hand, is faithfully copied from the original (excluding the parts explaining how to answer questions, for which we have adapted the instructions to our annotation interface).

## 4. Results

One of the main difficulties one faces in faithfully reproducing an experiment carried out by others lies in gathering all the necessary information. If they are not directly stated in the reference paper, it is necessary to seek clarification from the persons involved. However, during this exchange of infor-
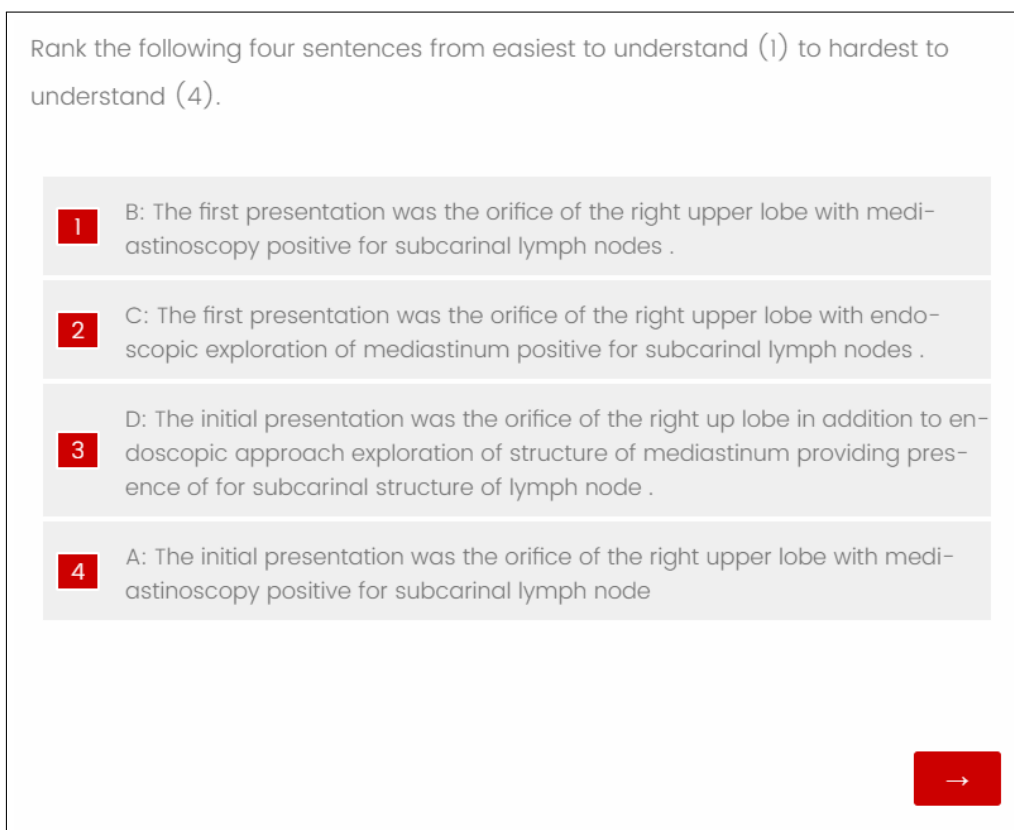
---

[3]https://www.prolific.com/

263

Rank the following four sentences from easiest to understand (1) to hardest to understand (4).

**1** B: The first presentation was the orifice of the right upper lobe with mediastinoscopy positive for subcarinal lymph nodes .

**2** C: The first presentation was the orifice of the right upper lobe with endoscopic exploration of mediastinum positive for subcarinal lymph nodes .

**3** D: The initial presentation was the orifice of the right up lobe in addition to endoscopic approach exploration of structure of mediastinum providing presence of for subcarinal structure of lymph node .

**4** A: The initial presentation was the orifice of the right upper lobe with mediastinoscopy positive for subcarinal lymph node

→

Figure 1: A screenshot of the annotation interface used in our replication study

| System | Rank:1 | | Rank:2 | | Rank:3 | | Rank:4 | | AVG | | CV* |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | O | R | O | R | O | R | O | R | O | R | |
| NTS + PT | 430 | 517 | 255 | 214 | 230 | 197 | 65 | 72 | **1.93** | **1.82** | 5.63 |
| NTS | 259 | 228 | 294 | 288 | 264 | 276 | 163 | 208 | 2.34 | 2.46 | 5.15 |
| ORIG | 120 | 123 | 222 | 233 | 381 | 408 | 257 | 236 | 2.79 | 2.76 | 1.19 |
| PTB | 171 | 132 | 209 | 265 | 105 | 119 | 495 | 484 | 2.94 | 2.96 | 0.51 |

Table 3: Comparison of original and reproduced results. *Rank:x* indicates the number of times each system was ranked at rank *x* and the last two columns show the average rank calculated according to the formula 1. O = Original and R = Reproduced. CV* is the Coefficient of Variation.

mation and material, doubts or misunderstandings may arise, as happened in this case: the results we are now going to present were, initially, completely different, due to a wrong assignment of the outputs to the 4 systems analysed. We find it interesting to mention this incident, as it is the consequence of one of the inherent difficulties of a reproducibility study such as this one.

**Side-by-side Results** Table 3 reports comparative results for the original (O) and reproduced (R) studies. It should be noted that the total number of annotations taken into account in the final results varies between the original experiment and our replication. This is due to the fact that, as mentioned in subsection 2.2, the authors of the original study had to remove 20 annotations, resulting in

980 final data points. In the reproduced results shown in Table 3, however, no annotations was removed (resulting in 1000 final data points), because all of them meet the response criteria.

What emerges from our study confirms the original results: the system proposed by the authors of the reference paper (NTS + PT) is the best performing one in their case, with an average rank of 1.93, and it is also the best one in our reproduction (1.82). Moreover, the general order of all systems turns out to be the same, with the Phrase Table Baseline as the worst performing one, generating outputs that are, in average, less understandable than the original sentences.

**Reproducibility Analysis** Following the protocol for the ReproHum project, in Table 3 we reported

| | Krippendorff's $\alpha$ | Pearson's $r$ | Spearman's $\rho$ |
|---|---|---|---|
| Agreement between Two Studies | 0.30 | - | - |
| IAA of Original Study | 0.22 | - | - |
| IAA of Replicated Study | 0.40 | - | - |
| Corr. between Two Studies (System Scores) | - | 0.98 | 1.00 |
| Corr. between Two Studies (Average Annotations) | - | 0.76 | 0.75 |

Table 4: Agreement between the two studies, calculated considering all 20 annotations for each sentence; IAA for the original and the replication study; correlation coefficients between the two experiments' results; correlation coefficients between the two experiments' sets of average annotations.

| System | Sentence | O | R |
|---|---|---|---|
| ORIG | A diagnostic paracentesis was said to show a sterile transudate. | 2.9 | 3.6 |
| NTS | A diagnostic paracentesis was said to show a good transudate. | 2.2 | 2.5 |
| NTS + PT | A diagnostic puncture and drainage was said to show a good transudate. | 1.3 | 1.1 |
| PTB | A diagnostic has intent puncture and drainage was said to show a sterile transudate. | 3.6 | 2.8 |
| ORIG | The tumor now involves the trachea as well as the right main bronchus. | 2.8 | 2.1 |
| NTS | The tumor now involves the opening as well as the right main bronchus. | 2.0 | 1.9 |
| NTS + PT | The tumor now involves the opening as well as the right main bronchial structure. | 1.5 | 2.3 |
| PTB | The tumor now involves the tracheal structure as good as the right main bronchial structure. | 3.7 | 3.7 |

Table 5: Examples of outputs produced by different systems and corresponding results from the original (O) and reproduced (R) rankings.

the Coefficient of Variation debiased for small sample size (CV*), as defined in Belz (2022).

We then calculated the agreement between ours and the original results, by considering all 20 annotators (10 from the original experiment and 10 from our reproduction study) for each sentence. We used the Krippendorff's $\alpha$ agreement measure as proposed in Castro (2017), and achieved an agreement of 0.30, as shown in Table 4. In the same table, we also reported the Inter-Annotator Agreement both within the evaluations collected by us and those collected by the authors of the reference study, for which we achieved higher scores.

We also calculated the correlation between the two sets of system final scores: ours and the original one, as reported in the column "AVG" of Table 3. Table 4 shows that a very high positive correlation was found, consistent with our similar results. To get more information on the quality of our reproduction, however, we also analysed the correlation between the two sets of single evaluations given by our annotators and the evaluations gathered in the original study. Specifically, we assigned each of the 400 annotated sentences (4 sentences for 100 instances) the average score received by the 10 annotators, and ran the correlation between the two studies. The results show that they have high correlation scores on both levels, confirming our results consistent with the original study. Lastly, we

reported an error count on these two lists of average rank, rounding the average rank to the nearest whole number, and found that 250 of the 400 values from the two studies agree, while 150 values differ.

**Case Study** Table 5 shows two examples of annotations: for each example, we reported the four evaluated outputs and the average score obtained by the ten annotators, both in the original experiment and in our reproduction. It can be seen that the NTS + PT system makes targeted changes to the original sentence, managing to modify too technical terms. In the first example, these changes result in increased understandability from the original sentence; however, for the second example, our annotators found the original sentence to be slightly more understandable. The baseline, on the other hand, makes a substantial number of changes, but these do not always help to increase the understandability of the sentence.

## 5. Conclusion

The main objective of this study was to remain as faithful as possible to the experimental choices made by the authors of the original paper when replicating the human evaluation they ran on system outputs. Any independent decisions we made were motivated by contingencies beyond our con-

trol (such as the use of a different crow-sourcing platform) or by a lack of information (e.g., concerning the compensation due to the annotators). Although the reproducing process is intrinsically difficult, the results we obtained align with the general findings of the original paper.

## Acknowledgments

## 6. Bibliographical References

Suha S. Al-Thanyyan and Aqil M. Azmi. 2021. Automated text simplification: A survey. *ACM Comput. Surv.*, 54(2).

Anya Belz. 2022. A metrological perspective on reproducibility in nlp. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz, Anastasia Shimorina, Shubham Agarwal, and Ehud Reiter. 2021. The ReproGen shared task on reproducibility of human evaluations in NLG: Overview and results. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 249–258, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Anya Belz, Anastasia Shimorina, Maja Popović, and Ehud Reiter. 2022. The 2022 ReproGen shared task on reproducibility of evaluations in NLG: Overview and results. In *Proceedings of the 15th International Conference on Natural Language Generation: Generation Challenges*, pages 43–51, Waterville, Maine, USA and virtual meeting. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, Gavin Abercrombie, Jose M. Alonso-Moral, Mohammad Arvan, Anouck Braggaar, Mark Cieliebak, Elizabeth Clark, Kees van Deemter, Tanvi Dinkar, Ondřej Dušek, Steffen Eger, Qixiang Fang, Mingqi Gao, Albert Gatt, Dimitra Gkatzia, Javier

González-Corbelle, Dirk Hovy, Manuela Hürlimann, Takumi Ito, John D. Kelleher, Filip Klubicka, Emiel Krahmer, Huiyuan Lai, Chris van der Lee, Yiru Li, Saad Mahamood, Margot Mieskes, Emiel van Miltenburg, Pablo Mosteiro, Malvina Nissim, Natalie Parde, Ondřej Plátek, Verena Rieser, Jie Ruan, Joel Tetreault, Antonio Toral, Xiaojun Wan, Leo Wanner, Lewis Watson, and Diyi Yang. 2023. Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP. In *Proceedings of the Fourth Workshop on Insights from Negative Results in NLP*, pages 1–10, Dubrovnik, Croatia. Association for Computational Linguistics.

Santiago Castro. 2017. Fast Krippendorff: Fast computation of Krippendorff's alpha agreement measure. https://github.com/pln-fing-udelar/fast-krippendorff.

Kevin Donnelly. 2006. Snomed-ct: The advanced terminology and coding system for ehealth. *Studies in health technology and informatics*, 121:279–90.

Alistair Johnson, Tom Pollard, Lu Shen, Li-wei Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Celi, and Roger Mark. 2016. Mimic-iii, a freely accessible critical care database. *Scientific Data*, 3:160035.

Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander Rush. 2017. OpenNMT: Open-source toolkit for neural machine translation. In *Proceedings of ACL 2017, System Demonstrations*, pages 67–72, Vancouver, Canada. Association for Computational Linguistics.

Sergiu Nisioi, Sanja Štajner, Simone Paolo Ponzetto, and Liviu P. Dinu. 2017. Exploring neural text simplification models. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 85–91, Vancouver, Canada. Association for Computational Linguistics.

Ricardo Rei, 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. Association for Computational Linguistics.

Matthew Shardlow and Raheel Nawaz. 2019. Neural text simplification of clinical letters with a domain specific phrase table. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 380–389, Florence, Italy. Association for Computational Linguistics.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Özlem Uzuner, Yuan Luo, and Peter Szolovits. 2007. Evaluating the State-of-the-Art in Automatic De-identification. *Journal of the American Medical Informatics Association*, 14(5):550–563.

## A   Appendix

Overview

In this task you must rank 4 sentences from easiest to understand to hardest to understand. The sentences are taken from clinical discharge letters and have been automatically processed to make them easier to understand. We want to find out which of the various methods we have used to improve the sentences is the best. You don't need to worry too much about small grammatical errors (i.e., punctuation in the wrong place, etc.), instead you should focus on the meaning and how well that meaning will be understood by a patient reading a letter sent home to them by their doctor. Typically this will be a case of judging whether the words that have been used are more likely to be understood by a patient without specialist medical expertise. This is a naturally subjective task and we expect you to use your own judgment to identify what would be easiest to understand for a patient reading this information in a letter from their doctor.

Steps

You will be presented with 4 sentences labelled A, B, C and D. You should first read the sentences carefully and ensure that you understand the meaning behind them. You will be asked to rank the sentences from easiest to understand to hardest to understand. You should put the sentence that you find the easiest to understand in the first line of the list. The next easiest goes in the second line, and so on. The sentence that you found the most difficult to understand should go in the fourth line. You must have a different sentence in each line. The differences in sentences may be small, but we still want you to make a judgement about which is better than the other. all four sentences should be different in every case, but if you find two sentences that are the same then just put them next to each other in the rankings, selecting the highest letter in the alphabet as the higher rank (i.e., A should be above B if and only if the sentences are completely identical).

Tip: take time to read the sentences and understand the meaning behind them.

Examples

A: The patient had a fractured tibia

B: The patient had a broken arm

C: The patient had a fractured arm

D: The patient sustained a fractured tibia

Ranking:

B is the easiest to understand for a patient (as it uses 'had', 'broken' and 'arm', which are more commonly understood words)

C is the next easiest to understand (it uses 'had' and 'arm', but also 'fractured' which may not be understood by a patient)

A is the third easiest, or second most difficult (it uses 'fractured tibia' which is hard to understand without medical expertise)

D is the hardest (it uses 'sustained' in place of 'had' which may be further confusing to the patient)

→

Figure 2: A screenshot of the instruction interface in our replication study.

# ReproHum #0087-01: A Reproduction Study of the Human Evaluation of the Coverage of Fact Checking Explanations

**Mingqi Gao, Jie Ruan, Xiaojun Wan**

Wangxuan Institute of Computer Technology, Peking University

{gaomingqi, wanxiaojun}@pku.edu.cn

ruanjie@stu.pku.edu.cn

## Abstract

We present a reproduction study of the human evaluation of the coverage of fact checking explanations conducted by Atanasova et al. (2020), as a team in Track B of ReproNLP 2024. The setup of our reproduction study is almost the same as the original study, with some necessary modifications to the evaluation guideline and annotation interface. Our reproduction achieves a higher IAA of 0.20 compared to the original study's 0.12, but discovers a mismatch between the IAA calculated by us with the raw annotation in the original study and the IAA reported in the original paper. Additionally, our reproduction results on the ranks of three types of explanations are drastically different from the original experiment, rendering that one important conclusion in the original paper cannot be confirmed at all. The case study illustrates that the annotators in the reproduction study may understand the quality criterion differently from the annotators in the original study.

**Keywords:** reproduction study, human evaluation, fact checking explanations

## 1. Introduction

These years have witnessed the concern about reproducibility issues in the field of NLP, especially human evaluation (Belz et al., 2023). In this paper, we present a reproduction study of human evaluation of the coverage of fact checking explanations (Atanasova et al., 2020), as a team in the Track B of ReproNLP Shared Task 2024 (Belz and Thomson, 2024).

The original study (Atanasova et al., 2020) formalizes fact checking as follows: Given a claim and some ruling comments, the model is required to predict the veracity label of the claim and also generation explanations. In the original experiments, human evaluation was performed to compare the quality of gold explanations and the explanations generated by two proposed models. The explanations were ranked by human annotators according to four quality criteria separately: Coverage, Non-Redundancy, Non-Contradiction, and Overall. After a discussion with the organizers of ReproNLP, we are asked to conduct a reproduction study only for Coverage.

## 2. Experimental Design

### 2.1. Original Experiment

LIAR-PLUS (Alhindi et al., 2018), a fact checking dataset based on PolitiFact [1], was used in the original study. Each instance of the dataset contains a claim, some ruling comments, a veracity label, an automatically extracted justification

as the gold explanation, and other metadata (e.g. speaker). There are six veracity labels: pants-fire, false, mostly false, half-true, mostly-true, and true.

The gold explanations in the dataset are abbreviated as **Just** in the original study. Besides, two explanation generation models are proposed: **Explain-MT** was trained jointly with veracity label prediction and **Explain-Extr** was trained separately.

**Selection of evaluation instances.** According to the original paper, 40 instances were randomly selected from the test set and three veracity explanations were collected for each of them. Each instance for human evaluation includes an instance ID, a claim, a veracity label, and three explanations. The ruling comments are excluded. Additionally, it is worth mentioning that after examining the original annotation interface (the Excel file), we find there are 80 instances included. Nevertheless, according to the raw annotation in the original experiment, only the first half was annotated by all three annotators.

**Participating annotators and compensation.** It is reported in the original paper that three annotators were involved but other information is not mentioned. According to the materials provided by the organizers of ReproNLP, none of the annotators were English native speakers. They were all colleagues of the authors and had previous experience with fact checking tasks. There is no information on whether and how much they were paid.

---

[1] https://www.politifact.com/

| id | claim | LABEL | justification 1 | justification 2 | justification 3 | Coverage | | | Non-redundancy | | | Non-contradictory | | | Overall | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2568.json | impleme | FALSE | sure that an | the Grand | lower the | 1 | 2 | 2 | 1 | 2 | 2 | 1 | 2 | 2 | 1 | 2 | 2 |
| 11923.json | will work | half-true | checking if | she will "work | Medium post, | | | | | | | | | | | | |
| 11025.json | thousan | FALSE | honest but | her state of | would have | | | | | | | | | | | | |
| 10085.json | on | FALSE | Mount | Mount | independent, | | | | | | | | | | | | |
| 9622.json | women | TRUE | ton of | ton of women | said that | | | | | | | | | | | | |
| 7834.json | in the | TRUE | matter, and | ng that the | the right to | | | | | | | | | | | | |
| 2205.json | 106,000 | TRUE | d health and | that we had | that we had | | | | | | | | | | | | |
| 8606.json | ans have | half-true | the health | the health | said, | | | | | | | | | | | | |
| 575.json | McCain | barely-true | 'intervening' | Airbus get | say two | | | | | | | | | | | | |

Figure 1: Annotation interface used in the original experiment. There are 80 instances in total and only the first ten are shown.

| id | claim | LABEL | justification 1 | justification 2 | justification 3 | Coverage | | |
|---|---|---|---|---|---|---|---|---|
| 2568.json | impleme | FALSE | sure that an | Grand Canyon | lower the | 1 | 2 | 2 |
| 11923.json | will work | half-true | checking if | she will "work | Medium post, | | | |
| 11025.json | thousand | FALSE | honest but | her state of | would have | | | |
| 10085.json | on Mount | FALSE | Mount | Mount | independent, | | | |
| 9622.json | women | TRUE | ton of women | ton of women | said that | | | |
| 7834.json | in the | TRUE | matter, and | g that the right | the right to | | | |
| 2205.json | 106,000 | TRUE | health and | we had over | we had over | | | |
| 8606.json | ans have | half-true | the health | the health | "Republicans | | | |
| 575.json | McCain | barely-true | 'intervening' is | Airbus get the | say two | | | |

Figure 2: Annotation interface used in our reproduction experiment. There are 40 instances in total and only the first ten are shown.

**Quality criterion.** The definition of the coverage of the explanation is as follows:

*Coverage*. The explanation contains important, salient information and doesn't miss any important points that contribute to the fact-check.

**Evaluation methods.** Given three different explanations (**Just**, **Explain-Extr**, and **Explain-MT**), the annotators were asked to rank 1,2,3 according to the criterion. It is noted in the evaluation guideline that if there is a tie and two explanations seem to have the same rank, the annotation should assign the same rank to them.

**Annotation interface.** The annotation was conducted through an Excel file, a screenshot of which is shown in Figure 1. In each row, the three explanations were randomized in terms of where they were placed to ensure fairness. Annotators were asked to record their ranks of the three explanations in the same row.

**Annotation procedure.** According to the information provided by the organizers of ReproNLP, there is no training process. Three participants were asked to read the evaluation guideline and then annotate the selected 40 instances separately.

**Inter-annotator agreement (IAA).** Krippendorff's $\alpha$ (Hayes and Krippendorff, 2007) was used to measure the IAA.

**Presentation of results.** For each type of explanation, the mean average ranks (MAR) by each annotator were presented. The average MAR of the three annotators was taken as the final result.

## 2.2. Reproduction Experiment

We were provided with an Excel file that included all the evaluation instances and an evaluation guideline. Both of them are exactly the same as the original experiment, which makes the setup of our reproduction experiment almost identical to the original experiment. The main differences from the original experiment are described below. For more details, please refer to the Human Evaluation Sheet (HEDS) (Shimorina and Belz, 2022) in supplementary materials [2].

**Modifications to the evaluation guideline and the annotation interface.** In the original study,

---

[2] They are also available at https://github.com/nlp-heds/repronlp2024.

| | Original | Reproduction | Confirmation |
|---|---|---|---|
| 1 | The gold explanation ranks the best in Coverage. | The gold explanation ranks the worst in Coverage. | Not confirmed. |
| 2 | *Explain-MT* ranks better than *Explain-Extr* in Coverage. | *Explain-MT* ranks better than *Explain-Extr* in Coverage. | Confirmed. |

Table 1: The conclusions from the original paper and the conclusions according to our reproduction results. The confirmation column shows whether the conclusion in the original study is confirmed or not.

| | Just | Explain-Extr | Explain-MT |
|---|---|---|---|
| original (calculated by us) vs. original (from the paper) | 1.34 | 1.60 | 0.59 |
| reproduction vs. original (from the paper) | 38.14 | 2.09 | 3.63 |
| reproduction vs. original (calculated by us) | 36.85 | 3.68 | 4.22 |

Table 2: CV*s among different experiment results. The smaller the CV*, the closer the results.

| | |
|---|---|
| original (calculated by us) vs. original (from the paper) | 1.00 |
| reproduction vs. original (from the paper) | -0.50 |
| reproduction vs. original (calculated by us) | -0.50 |

Table 3: Spearman's $\rho$ among different experiment results.

| | nominal | ordinal | interval | ratio |
|---|---|---|---|---|
| Original (calculated by us) | 0.05 | 0.12 | 0.12 | 0.12 |
| Reproduction | 0.12 | 0.20 | 0.20 | 0.18 |

Table 4: Krippendorff's $\alpha$. Different columns denote the annotations are viewed as nominal, ordinal, interval, or ratio data. In general, ranks are considered ordinal data.

in addition to Coverage, the annotators needed to assess the explanations against each of the three other quality criteria: Non-Redundancy, Non-Contradiction, and Overall. Additionally, there is another human evaluation task in the original study: providing the veracity label based on the explanations. These are reflected in the original evaluation guideline and the Excel file. We removed the content about other quality criteria and tasks from the evaluation guideline and the Excel file because we only reproduced the coverage evaluation of the explanations. The original evaluation guideline and the modified guideline are both included in the supplementary materials. The modified Excel sheet is shown in Figure 2. Furthermore, we only include the first 40 instances in our Excel file.

**Participating annotators and compensation.** Following the discussion with the organizers of ReproNLP, we recruited three PhD students who were proficient in English and paid them 12.24 EUR per hour.

## 3. Results

In addition to the evaluation guideline and the Excel file for annotation, we were also provided with

| Annotators | Just | Explain-Extr | Explain-MT |
|---|---|---|---|
| | Original (from the paper) | | |
| All | **1.48** | 1.89 | 1.68 |
| Annotator #1 | **1.50** | 2.08 | 1.87 |
| Annotator #2 | **1.74** | 2.16 | 1.84 |
| Annotator #3 | **1.21** | 1.42 | 1.34 |
| | Original (calculated by us) | | |
| All | **1.50** | 1.86 | 1.69 |
| Annotator #1 | **1.57** | 2.02 | 1.85 |
| Annotator #2 | **1.72** | 2.15 | 1.87 |
| Annotator #3 | **1.21** | 1.41 | 1.33 |
| | Reproduction | | |
| All | 2.18 | 1.93 | **1.62** |
| Annotator #1 | 2.05 | 2.00 | **1.65** |
| Annotator #2 | 2.23 | 1.95 | **1.62** |
| Annotator #3 | 2.25 | 1.85 | **1.60** |

Table 5: Mean Avegae Ranks (MAR) of different explanations, the lower the better. The best result of each row is in **bold** and the best generated explanation is in lightgrey . The first four rows are directly taken from the original paper. Note that there is no correspondence between the annotators of the original experiment and our reproduction experiment.

the raw annotation of each annotator in the original experiment, which enabled us to calculate the IAA and the MAR based on the raw annotation. It is interesting to see whether it is different from what is reported in the paper. Therefore, we present three types of experiment results in this section: the original results taken from the original paper, the original results calculated by us with the raw annotation, and our reproduction results. The former two share common human annotation, and the latter two share a common calculation process.

| Claim | Warren Buffett called (Jeff) Greenes scheme financial weapons of mass destruction. |
|---|---|
| Label | FALSE |
| Just | Still, there's no doubt Greene was heavily involved in buying credit default swaps, which are a type of derivative. And Buffett was exceptionally critical of derivatives. The problem with Meek's ad is that it gives the impression that Buffett criticized Greene personally, which he didn't do. By saying Buffett called Greene's scheme a bad idea, Meek personalizes a comment that had nothing to do with Greene. |
| Just MAR | Original: 1.00 (1, 1, 1), Reproduction: 3.00 (3, 3, 3) |
| Explain-Extr | Warren Buffett called Greene's scheme 'financial weapons of mass destruction.' We looked at Meek's statement that Greene "is the king of the undercover Credit Default Swap that brought about the destruction of our economy that we have right now." 317 on the Forbes list) would have made it onto Buffett's radar. We should note here that journalists love Buffett's phrase "financial weapons of mass destruction." |
| Explain-Extr MAR | Original: 2.67 (2, 3, 3), Reproduction: 1.33 (1, 1, 2) |

Table 6: An example that shows the different annotation results between the original experiment and our reproduction experiment. The rows of MAR list how the three annotators rank the explanation and the mean average ranks.

## 3.1. Inter-annotator Agreement

The Krippendorff's $\alpha$ reported in the original paper is 0.26. However, it does not match the results calculated by us with the raw annotation in the original experiment, as shown in Table 4. Based on our calculation, the IAA of our reproduction experiment is higher than the original experiment, but both of them are even lower than 0.26. After double-checking our calculation process through several computational scripts [3], we suspect that there may be some problems in how the IAA is calculated in the original study or there are some unknown details. Perhaps the annotation outside the first 40 instances was used.

The original paper considers a low IAA of 0.26 may be caused by the high subjectivity of ranking and the difficulty of this task. We believe that the inadequate evaluation guideline may also contribute to the low IAA. First, there is no example for each quality criterion. Second, the six veracity labels (pants-fire, barely-true, half-true, mostly-true, false, and true) lack clear definitions, which makes the evaluation of explanations harder.

## 3.2. Side-by-side Comparisons

Table 5 shows that there are minor differences in MAR between the results taken from the original paper and the results calculated by us with the raw annotation in the original experiment. However, **our reproduction results are dramatically dif-**

**ferent from the original experiments.** As shown in Table 1, a conclusion that the gold explanation ranks the best for Coverage is not confirmed at all, and our reproduction experiment yields the opposite conclusion. Despite this inconsistency, another conclusion is confirmed by our reproduction experiment.

We also present CV*, a metric proposed by Belz et al. (2022) to quantify reproducibility (in Table 2) and Spearman's $\rho$ (in Table 3) among different experiment results, also demonstrating the small differences between the original results calculated by us and from the paper but sharp inconsistency between our reproduction experiment and the original experiment.

## 3.3. Discussion

The big difference in the ranks of the gold explanations (Just) encourages us to conduct a case study. After examining some instances that differ from the original annotations, we conclude that the annotators in the reproduction study may understand the quality criterion differently from the annotators in the original study. The annotators in the original study pay more attention to whether the veracity label can be inferred from the explanation, while the annotators in the reproduction study focus more on whether the information in the claim is covered by the explanation. Table 6 shows an example. The annotators' understanding in the original study may be more reasonable but the ambiguity in the definition of the quality criterion is also the cause of this phenomenon.

---

[3]Our calculation results were verified by both a Python library (https://github.com/grrrr/krippendorff-alpha) and an online calculator (Marzi et al., 2024).

## 4. Conclusion

In this paper, we present a reproduction study of the human evaluation of the coverage of fact checking explanations under the guidance of the organizers of ReproNLP. Our conclusions are as follows:

- Our reproduction achieves a higher Krippendorff's $\alpha$ of 0.20 than the original experiment (0.12) based on our calculation, though both of them are not satisfactory.

- Krippendorff's $\alpha$ calculated by us with the raw annotation in the original experiment does not match what is reported in the original paper.

- The results of our reproduction experiment are drastically different from the original experiment, rendering that one important conclusion in the original paper cannot be confirmed at all.

- There are minor differences between the results calculated by us with the raw annotation in the original study and the results reported in the original paper.

## 5. Bibliographical References

Tariq Alhindi, Savvas Petridis, and Smaranda Muresan. 2018. Where is your evidence: Improving fact-checking by justification modeling. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, pages 85–90, Brussels, Belgium. Association for Computational Linguistics.

Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020. Generating fact checking explanations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7352–7364, Online. Association for Computational Linguistics.

Anya Belz, Maja Popovic, and Simon Mille. 2022. Quantified reproducibility assessment of NLP results. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16–28, Dublin, Ireland. Association for Computational Linguistics.

Anya Belz and Craig Thomson. 2024. The 2024 repronlp shared task on reproducibility of evaluations in nlp: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, and Simon Mille. 2023. Non-repeatable experiments and non-reproducible results: The reproducibility crisis in human evaluation in NLP. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3676–3687, Toronto, Canada. Association for Computational Linguistics.

Andrew F. Hayes and Klaus Krippendorff. 2007. Answering the call for a standard reliability measure for coding data. *Communication Methods and Measures*, 1(1):77–89.

Giacomo Marzi, Marco Balzano, and Davide Marchiori. 2024. K-alpha calculator–krippendorff's alpha calculator: A user-friendly tool for computing krippendorff's alpha inter-rater reliability coefficient. *MethodsX*, 12:102545.

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

# ReproHum #0866-04: Another Evaluation of Readers' Reactions to News Headlines

**Zola Mahlaza[1], Toky Raboanary[1], Kyle Seakgwa[1,2], C. Maria Keet[1]**

[1]University of Cape Town, Cape Town, South Africa
[2]University of the Western Cape, Cape Town, South Africa
{zmahlaza,traboanary,mkeet}@cs.uct.ac.za, SKGKYL001@myuct.ac.za

## Abstract

The reproduction of Natural Language Processing (NLP) studies is important in establishing their reliability. Nonetheless, many papers in NLP have never been reproduced. This paper presents a reproduction of Gabriel et al. (2022)'s work to establish the extent to which their findings, pertaining to the utility of large language models (T5 and GPT2) to automatically generate writer's intents when given headlines to curb misinformation, can be confirmed. Our results show no evidence to support two of their four findings and they partially support the rest of the original findings. Specifically, while we confirmed that all the models are judged to be capable of influencing readers' trust or distrust, there was a difference in T5's capability to reduce trust. Our results show that its generations are more likely to have greater influence in reducing trust while Gabriel et al. (2022) found more cases where they had no impact at all. In addition, most of the model generations are considered socially acceptable only if we relax the criteria for determining a majority to mean more than chance rather than the apparent $> 70\%$ of the original study. Overall, while they found that "machine-generated MRF implications alongside news headlines to readers can increase their trust in real news while decreasing their trust in misinformation", we found that they are more likely to decrease trust in both cases vs. having no impact at all.

**Keywords:** text generation, reproduction, misinformation

## 1. Introduction

The reproduction of Natural Language Processing (NLP) studies is critical in establishing the reliability of published findings. This is especially timely since there is evidence that a number of NLP studies are not repeatable (Belz et al., 2023). The low levels of replicability observed in these investigations warrants significant attention given that in many other fields, such as the social and medical sciences, low levels of replicability have also been observed in large scale replication efforts (OpenScienceCollaboration, 2015). The results in such studies triggered a decade-long reckoning with this "reproducibilty crisis" (Baker, 2016). It eventually led to more stringent standards being adopted for reporting results, a push toward preregistered research designs, and the adoption of more open science methods, like the sharing of datasets (Vazire, 2018).

While the low levels of replicability initially called into question the reliability of results in the social and medical sciences, efforts to address these shortcomings triggered what has been called a "credibility revolution" due to widespread adoption of the aforementioned improvements (Vazire, 2018). For the NLP community to undergo a similar "credibility revolution", more research like (Belz et al., 2023) needs to be undertaken to ascertain the extent of its reproducibility problem. As part of an effort to ascertain the extent to which existing work is reproducible (Belz and Thomson, 2024), this paper reports on the reproducibility of the human evaluation study conducted by Gabriel et al. (2022).

The work by Gabriel et al. (2022) focuses on investigating the utility of text generation models for automatically generating a writer's intent when given a news headline, as a means of combating misinformation. While the original work focuses on numerous tasks (e.g., it described the creation of a misinformation news headline corpus with human annotations of the writer's intent, readers' perception, possible actions that could be taken by the reader, and the likelihood of spread of the associated article), our sole focus in on the reproducibility of its human evaluations.

We investigate the reproducibility of the original study via a survey with 42 crowd-workers[1] who are based in the United States and judge the headline and intent pairs from the original study. The nature of the study is kept the same, where possible, and we compare the resulting findings to establish whether there is any difference with the original work. We have found that the results obtained with our survey contradict 2/4 of the findings from the original study and we can partially support two of the original study's findings. Specifically, with respect to the partially supported findings, most of the models' generations are considered socially acceptable if the criteria for determining a majority means more than chance[2] instead of $\geq 70$", a value

---

[1]One was excluded in the final analysis as they submitted incomplete survey responses

[2]We assume that "chance" means 50%

that can be inferred from the results. In addition, while all models were rated as being capable of influencing readers to trust or distrust, T5's generations are more likely to have greater influence in reducing trust while Gabriel et al. (2022) found more cases where they had no impact at all.

The rest of the paper is structured as follows. Section 2 summarises how the original study was conducted and lists the findings that emanated from it. Section 3 describes how the reproduction survey was set up, methods used to compare Gabriel et al. (2022)'s work with the current study, and the results we obtained in our survey. The differences and similarities with respect to findings between the two studies are discussed in Section 4, and Section 5 concludes.

## 2. Original study

While most work on combating misinformation pursues the creation of models to classify headlines, or articles, as being real or misinformation, Gabriel et al. (2022) takes a different approach towards building AI models. They investigate the extent to which machine-inferred writer's intents can improve reader's ability to identify misinformation. They do so by creating a human-annotated news corpus of headlines and intents with which they fine-tune pre-trained language models. The utility of the generated intents is evaluated by humans.

### 2.1. Dataset and models

The headlines were sourced from published misinformation datasets about Covid-19 (Cui and Lee, 2020; Gruppi et al., 2021; Network, 2024; Shapiro et al., 2020), climate change (Gruppi et al., 2021; Nørregaard et al., 2019), and cancer (Cui et al., 2020). The authors use the dataset to train models to automatically generate the writer's intent when given a headline and associated information (e.g., domain of the associated article/headline — either Covid-19, climate change, or cancer). The writers' intents are generated using two pre-trained language models, namely T5 (Raffel et al., 2020) and GPT2 (Radford et al., 2019).

### 2.2. Evaluation

The corpus and models are used to (1) investigate whether the headlines are trustworthy, be this with or without the writers' intent annotations that are automatically generated by the models, (2) determine whether the generated writers' intents are coherent and relevant, (3) establish whether the writers' intents are socially acceptable, and (4) ascertain whether the headlines and/or writers' intents perpetuate negative social biases or stereotypes.

## 2.3. Findings

Mechanical Turk (MTurk) workers' judgements of the trustworthiness of each news headline (with and without the intent) shows that while there were changes after seeing the intent, in the best case (i.e., intents generated by T5) there was only a weak positive correlation with the true class label (i.e., real/misinformation). Workers were also asked to judge the overall quality of the machine-generated intent in terms of coherence and its relevance to the headline on a 5-point Likert scale. The judgements show that the intents generated by T5 were perceived to have better quality, with an average of 3.74. Workers were also asked to judge whether the writer's intent conveys feelings or thoughts that are socially acceptable on a binary scale. In the case of one of the T5 variants, a model with the highest socially acceptable intents, we see a percentage of 75.30%. While worker's judgements of the capability of the beliefs and/or news events to perpetuate negative social biases or stereotypes were solicited, the results are not reported.

## 3. Reproducibility Study Design

The goal of our study was to reproduce the human evaluations using the same resources and methods as the original study, where possible. We did not aim to recreate their text generation models from scratch, but only reproduce the human evaluation thereof. We conducted the survey via Prolific[3] and an institutionally hosted version of LimeSurvey[4]. The human evaluation datasheet (Shimorina and Belz, 2022) for the study is shared via Github[5].

### 3.1. Survey

We created a survey using a dataset of 600 tuples of human authored headlines and automatically generated writers' intents. The dataset was sourced from Gabriel et al. (2022) via the organisers of the ReproNLP (Belz and Thomson, 2024) task. Each writer's intent is either 'real' or 'misinformation' and it is generated by one of the two types of models described in Section 2. Since the original study does not specify the number of texts evaluated by each participant, we split the dataset into 13 batches of 45 headline and intent pairs and one batch with 15 pairs. This was done to prevent collecting low quality judgements due to participant fatigue. Each batch was packaged into a survey where the participant is first given instructions, verbatim from the

---

[3] https://www.prolific.com/
[4] https://survey.cs.uct.ac.za/limesurvey/
[5] https://github.com/nlp-heds/repronlp2024

original study, describing what to expect as part of the survey ("You will read a sentence fragment describing a belief someone reading a news headline would have...") and what questions will be posed (e.g., "Please rate the quality of the belief description based on the following questions..."). They are then asked to judge quality of headline and intents, as shown in Figure 1.

Since the original study elicited 3 unique judgements per headline, we attempted to abide by the criteria as much as possible. We created a Python application (a web application created using the Django framework) to randomly assign a Prolific worker to one of the 14 surveys, provided it has less than three responses at the time of initiating the task. The survey was distributed to 42 Prolific participants that are based in the US, have 99% task approval, and have at least 200 tasks that have been approved.

**Evaluation strategy** There are two components to the evaluation. First, the calculations as by Gabriel et al. (2022). Overall Quality (coherence and relevance) is recorded on a 1-5 Likert scale. Influence on Trust is measured as more (+) or less (-) trustworthy, calculated as percentages, based on a 5-point scale that asked for the readers' perception. Third, for the sociopolitical acceptability, participants rate "their perception of the beliefs invoked by an implication in terms of whether they represent a majority (mainstream) or minority (fringe) viewpoint", where Gabriel et al "refer to "minority" viewpoint broadly in terms of less frequently adopted or extreme social beliefs, rather than in terms of viewpoints held by historically marginalized groups". This is reported as a percentage. We also recorded, on a nominal scale, the capacity of the headline and/or intent to perpetuate negative social biases or stereotypes. We report this as a percentage, even though it is not reported in the original study.

While the methods in Gabriel et al. (2022) do not describe further details, the results table indicates also "Corr w/ Label (all gens)", "Corr w/ Label (quality $\geq 3$)", and statistical significance. We take these to be correlations and a student-t test (with "$p < .05$").

Second, the comparisons of the results obtained in this reproduction are to be compared to the original results as reported in Gabriel et al. (2022). This involves both a numerical comparison and whether the same conclusions can be drawn from the results obtained. We first established where there is a difference in the computed percentages via a two-sample proportion hypothesis test (i.e., Z-test). We do not compare whether there is a significant difference between the means since they are only computed for Likert scales and are likely to lead to misinterpretations, especially since there is a

potential difference in how the evaluated data was batched.

Following that, we guided our comparisons using the findings (abbreviated **F** in the list) of the original study:

**F1:** "The T5-large model was rated as having slightly higher quality generations than the other model variants": We compared whether T5's average score was higher than the alternative model.

**F2:** "Most model generations were rated as being "socially acceptable"": We calculated whether most generations were judged as being socially acceptable. The original study does not specify what they deem a 'reasonable' majority is, so the cutoff point could belong to any value in the range (51-100), but we decided to use $\geq 70$.

**F3:** "All models were rated capable of influencing readers to trust or distrust": We counted the number of models for which the change in trust of distrust is non-zero.

**F4:** "[For] T5-base [, there is a] consistent correlation between the actual label and shifts in trustworthiness scores before and after annotators see the generated writer's intent": Working under the assumption that the original study established the consistency of a positive and significant correlation for T5-base between two binary attributes (shift in trust & label), we calculated the tetrachoric correlations[6] for the models and tested whether they are positive and significant[7] (two-tailed t-test, alpha=0.05) for T5-base but not the other models.

### 3.2. Results

The fourteen batches received a variable number of responses, as shown in Table 1. The results presented henceforth are not relevant to batch 10 since it did not get any responses. We also excluded a single response in batch 7 since one participant submitted incomplete responses.

Results of the participants' judgements, from the original and current study, are presented in Table 2. The results pertaining to the capacity of the intent and/or headline to perpetuate negative social biases or stereotypes are given in Table 3.

A comparison of the percentage of headlines that are considered socially acceptable, via the Z-test, shows that there is a significant difference (z = 2.51 for T5-base, z = 2.39 for T5-large, and z = 3.48

---

[6]Calculated using https://py-pair. readthedocs.io/

[7]H0: there is no linear relationship between the shift in trust and headline label.

Read a written description of a belief and rate the quality of the writing.

News Headline: Worshipers were beaten by policemen because they gathered in a church.

Belief description: The writer is implying that muslims are violent

**Does the belief description impact whether or not you think the news event is trustworthy?**

This is a question help text.

| | Very Trustworthy (Clearly real news) | Likely Trustworthy | Unsure | Likely Untrustworthy | Very Untrustworthy (Clearly misinformation) |
|---|---|---|---|---|---|
| Rate the trustworthiness of the headline before seeing the belief description: | ○ | ○ | ○ | ○ | ○ |
| Rate the trustworthiness of the headline after seeing the belief description: | ○ | ○ | ○ | ○ | ○ |

**What do you think the overall quality of the belief description is?**

(1 = Excellent, clearly coherent and relevant, 2 = Good, 3 = Fair, 4 = Poor, 5 = Nonsense. It doesn't make any sense given the event, is unreadable or is a copy of the headline)

○ 1    ○ 2    ○ 3    ○ 4    ○ 5

Figure 1: Screenshot of a section of the survey, with one of the news headlines, its belief description, and two of the questions for that headline and belief description.

Table 1: Number of responses for each batch. One of the responses from batch 7 was excluded because the crowdworker had not answered all the questions.

| Batch | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Responses** (#) | 8 | 2 | 2 | 2 | 3 | 2 | 5* | 2 | 2 | 0 | 3 | 2 | 3 | 5 |

for GPT2-large vs. critical value of 1.96) between the original study and our work for all three models. Specifically, the original study reported higher values of socially acceptable headlines. The extent to which the measured scores of social acceptability differ is not the same for all the models as Belz (2022)'s coefficient of correlation, given in Table 4, shows that GPT2-large exhibits the worst reproducibility while T5-large is better than the two alternative models.

Using the same test, we established that there is also a significant difference (z = -4.01 for T5-base, z = -3.80 for T5-large, and z = -1.09 for GPT2-large vs. critical value of 1.96) in the percentage of texts where there is an increase in trust after seeing the intent in the case of T5. However, we found no evidence that there is a significant difference in the case of GPT2. Noteworthy is that the original study recorded GPT2 as the model for whose intents have the greatest capacity to increase trust in the headline while the opposite was true in the current study (even if not statistically significant). The Z-test also showed that there is a significant difference (z = -40.09 for T5-base, z = -39.33 for T5-large, and z = -40.11 for GPT2-large vs. critical value of 1.96) in the percentage of texts for which there was

decrease in trust after seeing the intents.

The differences, with respect to a shift in trust, can be attributed to the high number of intents/headlines for which there was no change in participants' trust in the original study whereas no participant's trust was unaffected in our study.

We have found statistically significant evidence that there is no correlation between shift in trust and the class label in the case of T5-base, unless we exclude low quality (i.e., quality $< 3$) generations.

## 4. Discussion

We first compare our results to those reported in Gabriel et al.'s paper and then reflect on the reproducibility process.

### 4.1. Comparison of results with the original study

We now turn to confirm whether our study was able to confirm Gabriel et al. (2022)'s original four findings, as described in Section 3.1:

**F1:** Our results contradict this finding. Specifically, we found that GPT2, the alternative model, had higher quality generations than T5 (both

Table 2: Human evaluations from the original and reproduced study. Cells from original study marked with * indicate the statically significant existence of a correlation for $\alpha = 0.05$. Cell marked with ‡ indicate a statically significant the lack of a correlation for the same $\alpha$ value. Abbreviations: Orig = Original (i.e., Gabriel et al. (2022)), Corr = Correlation

| | Model | Quality (1-5) | Influence | | | | Socially accept. (%) |
| | | | +Trust (%) | -Trust (%) | Corr. (all gens) | Corr. (quality $\geq 3$) | |
|---|---|---|---|---|---|---|---|
| Orig. | T5-base | 3.61 | 8.33 | 7.82 | **0.24**$^*$ | **0.30**$^*$ | **75.30** |
| | T5-large | **3.74** | 7.73 | 9.76 | -0.03 | 0.09 | 74.66 |
| | GPT2-large | 3.46 | **9.70** | **13.10** | -0.04 | 0.10 | 74.66 |
| Ours | T5-base | 2.61 | **16.03** | 83.97 | 0.07$^‡$ | **0.99** | **68.67** |
| | T5-large | 2.56 | 14.77 | 85.43 | **0.99** | **0.99** | 68.31 |
| | GPT2-large | **2.77** | 11.68 | **89.38** | **0.99** | **0.99** | 65.30 |

Table 3: Percentage of headlines and intents that perpetuate negative stereotypes. Abbreviations: Sent = Sentence

| Model | Both do | Neither do | Sent. | News event |
|---|---|---|---|---|
| T5-base | 13.27 | 74.16 | 5.84 | 6.73 |
| T5-large | 18.4 | 73.09 | 7.96 | 5.49 |
| GPT2-large | 14.51 | 68.85 | 11.5 | 5.13 |

Table 4: Precision results for the socially acceptable attribute between the original and current study. Abbreviations: Unb. stdev = unbiased standard deviation, CV* = Belz (2022)'s coefficient of variation

| Model | Mean | Unb. stdev | CV* |
|---|---|---|---|
| T5-base | 71.985 | 5.876 | 9.1827 |
| T5-large | 71.485 | 5.628 | 8.8564 |
| GPT2-large | 69.980 | 8.295 | 13.3352 |

T5-base and T5-large). GPT2 also had the highest CV*, denoting the poor reproducibility when compared to the alternative models. In addition, the quality judgements of all the models were lower by about 1 point on a 5-point Likert scale.

**F2:** It is not clear what value is used by Gabriel et al. (2022) to determine a majority and our interpretation of their results suggests that they used $\geq 70$. Based on that interpretation, our results do not support this finding as we had fewer generations that are socially acceptable. There is a difference of 6-9 percentage points between our results and theirs and the difference is statistically significant. This finding

can only be supported to relax the cut-off point from $\geq 70$ to $\geq 51$.

**F3:** Our results confirm that all the models are judged to be capable of influencing readers' trust or distrust. However, there was a significant difference in T5's capability to reduce trust. Specifically, our results show that its generations are more likely to have greater influence in reducing trust while Gabriel et al. (2022) found more cases where they had no impact at all. There were also more cases where T5-base's generations positively shifted trust vs. GPT2.

**F4:** Our results contradict this finding. In fact, they show that T5-base is the only model for which there is no consistent correlation between the actual label and shifts in trustworthiness scores. When low quality (i.e., quality $< 3$) generations are included, our evidence demonstrates that there is no linear relation between the attributes. T5-large and GPT2 are the only models for which there is a consistent and strong correlation.

## 4.2. Challenges reproducing the study

The experiment's methods could not be reproduced exactly, due to several reasons. First, the ReproNLP project (Belz and Thomson, 2024) moved to Prolific (cf. MTurk in the original study). This had consequences for technically setting up the task. We cannot directly use the form created and used by Gabriel et al. (2022) for MTurk in Prolific. We had to use experimental software compatible with Prolific, such as LimeSurvey and Gorilla, or develop a new form hosted on another server and link it to Prolific to create a similar survey. We decided to use LimeSurvey to facilitate the experiment while maintaining the same objective. This allowed us to rely on its existing features such as recording

responses and collection of metadata pertaining to submission times, time taken to complete the survey, etc. However, the layout differs from Gabriel et al. (2022)'s survey.

Second, it was our decision to break up the task into batches, which may, or may not, have been done in the original study, as described in Section 3. With an estimated task completion time of 45 minutes for a batch of 45 headlines, it was deemed unreasonable to make a participant assess a batch of 600 headlines since that would have taken approximately 10 hours.

The change in evaluation platform used also resulted in a difference in the number of participants who evaluated each headline/intent pair per survey between the two studies. While our evaluation instrument was set up to abide by the upper limit of 3 responses per survey as much as possible, as mentioned in Gabriel et al. (2022), via assigning each worker to a batch that did not have $\geq 3$ responses already. We still obtained more than three responses for some batches, as included in Table 1, since there were cases where some workers were assigned to batches for which there were other users who were already evaluating but had not submitted their responses.

It is also possible that the way the headline and intent are shown to each participant may differ from the original study. We had access to a screenshot of the instrument used by Gabriel et al. (2022) and we determined that it uses a template (i.e., "News Headline: $[sentence]") at the top of the survey to display the information to be evaluated. However, it was unclear how and when the intent was presented to each participant. As such, we took the decision to include the intent alongside the news headline, as shown in Figure 1.

The instrument used by Gabriel et al. (2022) solicited judgements on a nominal scale, to determine whether the writer's intent and/or headline perpetuates negative social biases or stereotypes. However, the solicitation of those judgements was not described in the publication nor was there a presentation of the associated results. Nonetheless, we decided to collect and report such judgements for completeness.

An average cost of 12.60GBP (max of 12.62GBP) was spent on renumerating each participant and the figure includes the Prolific service fee of 3GBP and a value added tax between 0.59-0.61GBP. This was in line with the ReproNLP task's (Belz and Thomson, 2024) instructions which mandated a value of 12GBP per hour. The extent to which this figure differs from the original study is unclear since the original study does not specify how much evaluators were compensated. It only mentions the workers were paid $.6 per human intelligence task. However, it is unclear how many items were evaluated by each participant. In addition, of the 600 items to be evaluated, the authors excluded 12 since they were deemed as malformed or unsuitable, but the exclusion criteria were not reported/communicated.

Gabriel et al. (2022) note that they "obtained an Institutional Review Board (IRB) exemption for annotation work, and ensured annotators were fairly paid given time estimations.", however it is unclear if ethics approval was obtained for the evaluation task. It is possible that the authors either use the same term "annotator" to refer to both annotators and evaluators, when discussing ethics approval, since the latter are a subset of the crowd-workers that were initially recruited, or they may have deemed it unnecessary to seek ethics approval for the evaluation. The University of Cape Town does require ethics approval for experiments involving humans or large datasets, however, and thus needed to be obtained from the Science Faculty Ethics Committee. Besides filling in the form, this involved writing from scratch a research proposal, a data management plan, and a task-adapted consent form. It was approved with reference number SCI/00635/2024.

## 5. Conclusion

Our reproduction of the human evaluation component of the research reported in Gabriel et al. (2022) contradicts two out of the four findings reported in Gabriel et al. (2022), being **F1** and **F4**. Specifically, GPT2's generations had better quality even though Gabriel et al. (2022) found that T5's generations are better (**F1**) and T5 did *not* exhibit consistent correlations between the actual label and shifts in trustworthiness scores due to the inclusion of low quality (i.e., quality $< 3$) intents (**F4**). Most model generations were not rated as being "socially acceptable" (**F2**), unless one were to lower the cutoff point for determining a majority from $\geq 70\%$ to $\geq 51\%$. Lastly, our results confirm that all the models were capable of influencing readers' trust or distrust (**F3**). However, unlike the original study, T5's generations are more likely to have greater influence in reducing trust even though Gabriel et al. (2022) found more cases where they had no impact at all and there were more cases where T5-base's generations positively shifted trust vs. GPT2.

## 6. Acknowledgements

## 7. Bibliographical References

Monya Baker. 2016. 1,500 scientists lift the lid on reproducibility. *Nature*, 533:7604:452—4:454.

Anya Belz. 2022. A Metrological Perspective on Reproducibility in NLP*. *Computational Linguistics*, 48(4):1125–1135.

Anya Belz and Craig Thomson. 2024. The 2024 ReproNLP shared task on reproducibility of evaluations in NLP: Overview and results. In *Proceedings of the 4th Workshop on Human Evaluation of NLP Systems*.

Anya Belz, Craig Thomson, Ehud Reiter, and Simon Mille. 2023. Non-repeatable experiments and non-reproducible results: The reproducibility crisis in human evaluation in NLP. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 3676–3687. Association for Computational Linguistics.

Limeng Cui and Dongwon Lee. 2020. Coaid: Covid-19 healthcare misinformation dataset.

Limeng Cui, Haeseung Seo, Maryam Tabar, Fenglong Ma, Suhang Wang, and Dongwon Lee. 2020. Deterrent: Knowledge guided graph attention network for detecting healthcare misinformation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '20, page 492–502, New York, NY, USA. Association for Computing Machinery.

Saadia Gabriel, Skyler Hallinan, Maarten Sap, Pemi Nguyen, Franziska Roesner, Eunsol Choi, and Yejin Choi. 2022. Misinfo reaction frames: Reasoning about readers' reactions to news headlines. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 3108–3127. Association for Computational Linguistics.

Maurício Gruppi, Benjamin D. Horne, and Sibel Adalı. 2021. Nela-gt-2020: A large multi-labelled news dataset for the study of misinformation in news articles.

International Fact-Checking Network. 2024. Fighting the Infodemic: The CoronaVirusFacts Alliance. https://www.poynter.org/coronavirusfactsalliance/. [Online; accessed 01-April-2024].

Jeppe Nørregaard, Benjamin D. Horne, and Sibel Adalı. 2019. Nela-gt-2018: A large multi-labelled news dataset for the study of misinformation in news articles. *Proceedings of the International*

AAAI Conference on Web and Social Media, 13(01):630–638.

OpenScienceCollaboration. 2015. Estimating the reproducibility of psychological science. *Science*, 349:6251:943–6251:951.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

Colin Raffel, 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. *J. Mach. Learn. Res.*, 21:140:1–140:67.

Jacob N. Shapiro, Jan Oledan, and Samikshya Siwakoti. 2020. Fighting the Infodemic: The CoronaVirusFacts Alliance. https://esoc.princeton.edu/publications/esoc-covid-19-misinformation-dataset. [Online; accessed 01-April-2024].

Anastasia Shimorina and Anya Belz. 2022. The human evaluation datasheet: A template for recording details of human evaluation experiments in NLP. In *Proceedings of the 2nd Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 54–75, Dublin, Ireland. Association for Computational Linguistics.

Simine Vazire. 2018. Implications of the credibility revolution for productivity, creativity, and progress. *Perspectives on Psychological Science*, 13:4:411—4:417.

# Author Index