

Quality-aware Neural Machine Translation with Self-evaluation

Jiajia Cui¹, Lingling Mu¹, Qiuhui Liu², Hongfei Xu^{1†}

¹ Zhengzhou University, Henan 450001, China

² China Mobile Online Services, Henan 450001, China

jjcui@gs.zzu.edu.cn, iellmu@zzu.edu.cn

{liuqhano, hfxunlp}@foxmail.com

Abstract

The performance of neural machine translation relies on a large amount of data, but crawled sentence pairs are of different quality. The low-quality sentence pairs may provide helpful translation knowledge but also teach the model to generate low-quality translations. Making the model aware of the quality of training instances may help the model distinguish between good and bad translations while leveraging the translation knowledge. In this paper, we evaluate the quality of training instances with the average per-token loss (negative log-likelihood) from translation models, convert the quality scores into embeddings through vector interpolation and feed the quality embedding into the translation model during its training. We ask the model to decode with the best quality score to generate good translations during inference. Experiments on the IWSLT 14 German to English, WMT 14 English to German and WMT 22 English to Japanese translation tasks show that our method can effectively lead to consistent and significant improvements across multiple metrics.

Keywords: Neural Machine Translation, Quality-aware modeling

1 Introduction

The Transformer translation model (Vaswani et al., 2017) can produce high quality translations with large amounts of parallel data. Most large-scale parallel corpora are automatically extracted from crawled parallel websites and cannot ensure the quality of translation pairs. Training Neural Machine Translation (NMT) models on the unfiltered ParaCrawl leads to degraded translation quality (Junczys-Dowmunt, 2018; Schamper et al., 2018).

Peter et al. (2023) improve the NMT performance by filtering half of the training set with neural Quality Estimation (QE) metrics. However, this approach relies on the availability of high-performance QE models of the target language pair and fully discards the translation knowledge within the low-quality parallel data. Concurrently, Tomani et al. (2024) divide the training set into a number of bins based on the MetricX-QE scores and prompt the model with quality tags of the corresponding bins. However, 1) the quality embedding for each bin is only trained on part of the training set, 2) the quality differences inside a bin are fully neglected, 3) close quality scores divided into adjacent bins may have different quality embeddings, and 4) the method also relies on QE tools which may not be available for some language pairs.

To avoid the use of QE tools, we train NMT models on the parallel data for both directions, and use the bi-directional average per-token loss from the translation models as the quality measurements of training instances. We derive the quality embedding through vector interpolation with the quality score, and replace the embedding of the special start-of-sentence (< sos >) token with the quality embedding to make the model aware of the quality of training instances. All training data and corresponding translation

[†]Corresponding author: Hongfei Xu.

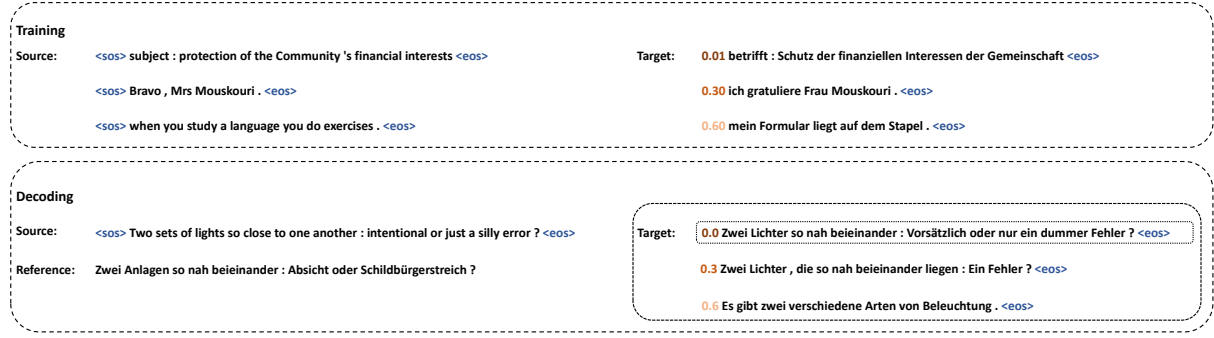


Figure 1: Quality-aware training and decoding. The left and right sides are for encoder and decoder respectively. The quality score is the first input of the decoder to trigger the auto-regressive decoding. Quality scores are converted into quality embeddings through linear interpolation. The best quality score is used for evaluation. Examples decoded with worse quality scores are only to show that the model follows the quality scores during translation.

knowledge are thus kept during the training. During inference, we prompt the model with the best quality score, and ask the model to generate good translations. Our main contributions are as follows:

- We train NMT models on the parallel data, and evaluate the data quality based on bi-directional scores to avoid the use of QE tools.
- We turn quality scores into quality embeddings by vector interpolation, and integrate into the translation model to help distinguish the quality of training instances while leveraging the translation knowledge.
- Our method brings consistent and significant improvements from low-resource to high-resource tasks across several evaluation metrics without negative impacts on the inference speed.

2 Our Method

The training and decoding of quality-aware translation are shown in Figure 1. In general, we first train source-to-target and target-to-source vanilla translation models (baseline without quality scores). Then we use the baseline models to calculate the average per-token prediction loss of each training instance as quality scores. The quality scores are used for the training of the quality-aware model by replacing the `<sos>` embedding with the quality embedding generated based on the quality score.

2.1 Bi-directional Self-evaluation

Fomicheva et al. (2020) show that NMT models are also strong quality estimators and can achieve good correlation with human quality judgments.

We use the baseline NMT models to compute the average per-token loss of each training instance. The prediction loss of the target token y_n is the negative log-likelihood computed by the NMT model \vec{M} based on the source sentence x and all preceding tokens $y_{<n}$ in the reference translation.

$$\text{loss}(x, \vec{M}, y_n) = -\log P(y_n | x, y_{<n}, \vec{M}) \quad (1)$$

But the loss of the forward direction may not fully reflect translation quality, especially when the source sentence is inadequately translated. While the target-to-source NMT model \overleftarrow{M} is likely to produce high loss scores for under-translated source tokens given the target sentence as encoder input. The bi-directional loss (loss_{bi}) considers both directions and is computed with Equation 2.

$$\text{loss}_{bi} = \frac{\sum_{t=1}^{|y|} \text{loss}(x, \vec{M}, y_t) + \sum_{t=1}^{|x|} \text{loss}(y, \overleftarrow{M}, x_t)}{|y| + |x|} \quad (2)$$

where $|x|$ means the number of reference tokens in the source sentence x .

We normalize the bi-directional average per-token loss into the range of $[0, 1]$ as quality score q with Equation 3.

$$q = \frac{loss_{bi} - \min(loss_{bi})}{\max(loss_{bi}) - \min(loss_{bi})} \quad (3)$$

There is no existing theoretical basis between the amount of data and the estimation reliability to our knowledge. But even the low-resource machine translation task (e.g., IWSLT 14 De→En) has a much larger training set (174k sentence pairs) than the amount of available training data for the quality estimation task and many other NLP tasks. Our method empirically works well on the low-resource IWSLT 14 De→En task as shown in Table 1.

2.2 Quality-aware NMT

The NMT model represents discrete tokens with embeddings, but the quality score for each instance is a continuous scalar. To make the NMT model explicitly aware of the training instance quality, we have to turn the quality score q into an embedding and feed the quality embedding into the model.

As the quality score is normalized into a fixed range ($[0, 1]$), we generate quality embeddings by vector interpolation. Specifically, we employ two vectors v_e and v_s as corresponding embeddings for the quality scores of 0 and 1 respectively. The embedding e_q of the other quality score q are computed by the weighted combination of v_s and v_e , where the weights are the distances from q to 0 and 1 respectively, as shown in Equation 4.

$$e_q = q * v_s + (1 - q) * v_e \quad (4)$$

As the auto-regressive decoder computes in a left-to-right manner, we replace the special $\langle \text{sos} \rangle$ embedding with the quality embedding to provide the model with quality information. We do not append the quality scores as additional word embedding dimensions to avoid changing the embedding size or reducing the number of dimensions for word representation. We also do not insert one more token before the $\langle \text{sos} \rangle$ token which may increase the complexity of decoding when applying the method to the decoder.

During inference, we feed the quality score representing the best quality into the model to generate good translations.

3 Experiments

3.1 Settings

Datasets. IWSLT 14 (Cettolo et al., 2014) using TED talks is the cleanest. Common Crawl is the major part of WMT 14 English (En) to German (De) task (Bojar et al., 2014) and noisy (COMET scores of 11.23% and 1.7% of the training set are below 0.5 and 0.3 respectively). ParaCrawl constitutes the largest part of the training set of the WMT 22 English to Japanese (Ja) task (Kocmi et al., 2022). Larger training sets normally tend to be more noisy. We tokenized and truecased English and German with Moses, and tokenized Japanese with MeCab (KUDO, 2010). We applied joint Byte-Pair Encoding (BPE) (Sennrich et al., 2016b) with 16k merge operations for the low-resource IWSLT 14 De→En task (174k), 32k merge operations for the WMT 14 En→De task (4.5M), and independent BPE with 32k merge operations for the WMT 22 En→Ja task (33M).

Baselines. We compared our method with the Transformer baseline, QE-based filter (Peter et al., 2023) and quality-aware prompting method (Tomani et al., 2024). For the Filter method (Peter et al., 2023), we filtered half of the training set with lower COMET-QE scores. For the Prompt method (Tomani et al., 2024), we used COMET-QE scores to divide the training set into 10 bins and replaced the $\langle \text{sos} \rangle$ token of both source and target sentences with bin tags following their paper. Recurrent decoders may lead to improved translation quality compared to the Transformer decoder (Chen et al., 2018). We also

tested the performance of our approach with the MHPLSTM decoder (Xu et al., 2021). Specifically, We replace the self-attention layers of the Transformer decoder with the MHPLSTM.

Hyperparameters. We followed the Base setting (512/2048, 8 heads) of Transformer (Vaswani et al., 2017) for WMT 14 En→De and WMT 22 En→Ja tasks. For IWSLT De→En task, we followed the experiment settings of Araabi and Monz (2020). For the training, we trained for 100k steps with a batch size of around 36k target tokens for the WMT 14 En→De task and the En→Ja task, a batch size of around 6k target tokens for the IWSLT 14 De→En task. The results in Table 1 are reproduced under the same setting. All the experiments used the same training script and hyper-parameters.

Evaluation. we used a beam size of 4 for decoding with the averaged model of the last 5 checkpoints saved with an interval of 1500 training steps. We evaluated with BLEU and chrF implemented by the sacreBLEU toolkit (Post, 2018) and COMET score (Rei et al., 2022a). Human evaluation would be valuable, but it is often impractical due to inconsistent standards across studies and neural metrics like COMET have led to high correlation with humans. For our method, we prompt the best quality score for high-quality translations during decoding and evaluation. The best quality score depends on the metric, it is 0 for bi-directional translation loss and MetricX, but 1 for COMET.

We found that the performance of Transformer is quite sensitive to small variations on the quality scores. It can vary for around 5 BLEU on the development and test sets when gradually increasing the decoding quality score from 0 to 0.1. So we treated the quality scores (the concatenation of q and $1 - q$ in Equation 4) like hidden representations in the neuron network and applied dropout to them during training to prevent the model from over-fitting the quality scores.¹

3.2 Main Results

We used bi-directional NMT losses as quality scores (§ 3.5) and added quality scores to the decoder (§ 3.6) based on our ablation study.

Results in Table 1 show that: 1) data filtering or prompting based on quality scores can lead to consistent and significant improvements across all metrics only on the high-resource WMT 22 En→Ja task of the largest training set, but their performances are close to the baseline on the WMT 14 En→De task and worse on the low-resource IWSLT 14 De→En task with some metrics, while our method is also effective for low-resource and middle-resource tasks, obtaining consistent and significant improvements with all metrics, and 2) our method can also improve the performance of the stronger recurrent decoder baseline.

Methods	IWSLT 14 De→En			WMT 14 En→De			WMT 22 En→Ja		
	BLEU	chrF	COMET	BLEU	chrF	COMET	BLEU	chrF	COMET
Transformer	29.57	52.45	77.67	27.64	57.26	82.42	21.55	29.47	83.76
Filter (2023)	28.39	51.51	75.66	27.68	57.21	82.75	21.92	30.02	84.31
Prompt (2024)	29.45	52.66	77.09	27.92	57.43	82.74	22.06	30.05	84.34
Ours	29.99	52.85	78.01	28.78	57.91	83.51	22.56	30.66	84.83
MHPLSTM (2021)	29.90	53.20	78.19	28.36	57.57	83.01	22.12	30.10	84.53
Ours	30.34	53.42	79.19	29.01	57.94	83.55	22.75	30.69	84.88

Table 1: Main results. Bold scores indicate improvements of our method over corresponding baselines.

3.3 Results on Large Language Model

We tested the performance of our method with the Large Language Model (LLM) on the WMT 22 En→Ja task, which is the largest dataset. We fine-tuned Qwen3-8B (Team, 2025) with LoRA with rank $r = 128$ as our baseline. LoRA was applied to the attention query/key/value projection layer. Experiments were conducted with identical configurations, including the number of training steps and batch size. We incorporated quality scores into the LLM fine-tuning for our method. Specifically, we

¹The model leads to BLEU scores of 23.72, 28.43 and 28.36 on the WMT 14 En→De task when decoding with 0, 0.05 and 0.1 as the quality scores respectively without dropout, while obtaining a highest BLEU score of 28.78 with a quality score of 0 with dropout probability of 0.1.

replaced the embedding of the start tag ($\langle im_start \rangle$) that triggers the generation with the quality embedding. We used a beam size of 4 for decoding with the best performing model on the development set.

Results in Table 2 show that our method outperforms the baseline across all metrics (+1.25/+1.48/+1.24 on BLEU/chrF/COMET). The BLEU and chrF scores are lower than the results in Table 1, but the COMET score is higher, which suggests that LLM generates more semantically adequate translations compared to the Transformer.

Method	BLEU	chrF	COMET
Baseline	20.49	28.28	85.85
Ours	21.74	29.76	87.09

Table 2: LLM Result on WMT 22 test set.

3.4 Effectiveness with Scaling

We tested the performance of shallower and deeper models on the WMT 14 En→De task. Results in Table 3 show that our method can also obtain consistent improvements with varying depths, and shallow models with our method can perform comparably to deeper baselines.

Depth		Method	BLEU	chrF	COMET
Encoder	Decoder				
3	3	Baseline	26.52	56.15	80.58
		Ours	27.36	56.92	81.98
6	6	Baseline	27.64	57.26	82.42
		Ours	28.78	57.91	83.51
12	6	Baseline	28.63	57.81	83.36
		Ours	29.48	58.55	84.37
18	6	Baseline	29.06	58.16	83.50
		Ours	29.84	58.67	84.39

Table 3: Results with varying depths.

3.5 Effects of Quality Estimation Methods

To verify that the translation losses (of the forward, reverse and bi-direction) are valid quality indicators, we measure the Pearson correlation coefficients between the translation losses and the COMET (Rei et al., 2022b)/MetricX (Juraska et al., 2023) scores. We compare with automatic metrics, because our method only uses quality scores of the training set. The vanilla translation models are trained without quality annotations and cannot overfit to quality scores. The quality estimation performance on the other datasets may not be robust due to the machine translation training data limitation, but our method only requires sufficient quality estimation performance on the training set.

Results in Table 4 show that: 1) translation losses of individual directions lead to similar Pearson coefficients to the COMET scores and the bi-directional translation loss leads to a slightly higher coefficient, and 2) the Pearson coefficient of the bi-directional translation loss (0.80) is close to that of MetricX (0.83), supporting that the translation losses can be used as effective quality measures. The correlation with MetricX scores is lower than with COMET for MetricX scores are somehow problematic (17.9% are 0, 8% are 25).

We conducted experiments on the WMT 14 En→De task to test the effects of different quality scores. Results in Table 5 show that: 1) all quality evaluation methods can lead to consistent improvements except for the MetricX, the poor performance with MetricX may be because more than 1/4 of the training set shares the same quality embeddings (17.9% and 8% of MetricX scores are 0 and 25 respectively) and

Target	Forward	Reverse	Bi-direction	MetricX
COMET	0.77	0.77	0.80	0.83
MetricX	0.70	0.71	0.73	-

Table 4: Pearson correlation coefficients between translation losses, MetricX and COMET/MetricX scores.

lacks discrimination in quality, 2) using the bi-directional loss is better than using a single direction loss, and 3) using COMET-QE scores leads to slightly higher COMET score but slightly worse BLEU and chrF scores than the bi-directional translation loss. This may be due to potential bias between COMET-QE and COMET. Using COMET-QE scores can obtain comparable performance to bi-directional loss and avoid the additional training cost of the bi-directional translation models for quality estimation, but self-evaluation with the translation loss can avoid the reliance on the existence of QE tools.

Methods	Dev Set			Test Set		
	BLEU	chrF	COMET	BLEU	chrF	COMET
Baseline	26.13	54.32	82.24	27.64	57.26	82.42
MetricX	25.69	53.93	81.79	27.97	57.17	82.26
COMET-QE	26.31	54.69	83.64	28.63	57.84	83.57
Forward	26.22	54.44	83.15	28.45	57.64	83.37
Reverse	26.33	54.64	83.31	28.50	57.82	83.35
Bi-direction	26.43	54.76	83.48	28.78	57.91	83.51

Table 5: Ablation study of quality estimation methods on the development set (newstest 2013) and test set (newstest 2014).

3.6 Quality-aware Encoding and Decoding

We tested the effects of providing quality scores to the encoder or the decoder or both on the WMT 14 En→De task. Results in Table 6 show that: 1) providing quality scores to either encoder or decoder can lead to significant improvements than the baseline, showing the importance of integrating quality scores, and 2) using quality scores in the decoder leads to better performance than in both, probably because that it might be better to only let the decoder learn to generate translations of quality than additionally asking the encoder to perform source language understanding for decoding of varying quality.

Methods	Dev Set			Test Set		
	BLEU	chrF	COMET	BLEU	chrF	COMET
Baseline	26.13	54.32	82.24	27.64	57.26	82.42
Encoder	26.18	54.60	83.04	28.41	57.54	83.19
Decoder	26.43	54.76	83.48	28.78	57.91	83.51
Both	26.21	54.60	83.08	28.54	57.73	83.13

Table 6: Results for quality-aware encoding/decoding on the development set (newstest 2013) and test set (newstest 2014).

3.7 Effects of Quality Scores on Inference

To test whether the model is aware of the translation quality after training, we tested the BLEU score of model translations with varying quality scores on the WMT 14 En→De task.

Results in Figure 2 show that the model learns to generate translations of the given quality score.

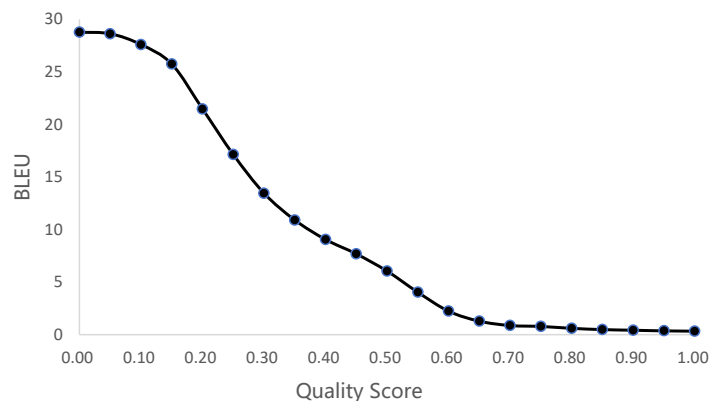


Figure 2: Translation quality with increasing loss scores.

4 Related Work

The performance of machine translation models heavily depends on the quality of the training data. Traditional approaches often improve data quality by filtering noisy data (Junczys-Dowmunt, 2018; Schamper et al., 2018; Chaudhary et al., 2019). Recently, Peter et al. (2023) use QE metrics to filter low-quality sentence pairs in the training data. And they found QE methods focus on selecting the best translation examples and identify more fine grained problems in the training data.

On the other hand, some studies have attempted to incorporate quality signals into the model training or decoding process. Fernandes et al. (2022) using QE for reranking or Minimum Bayes Risk decoding lead to significant improvements over beam search decoding. Concurrently, Tomani et al. (2024) eliminate the need for an external QE model during decoding by embedding quality signals directly into the model, they divide the training set into a number of bins based on the quality estimation scores and prompt the model with the quality embedding of the corresponding bin.

The evaluation methods for machine translation can be divided into two categories: reference-based and reference-free. Reference-based evaluation (BLEU (Papineni et al., 2002), BLEURT (Sellam et al., 2020) and COMET (Rei et al., 2022a)) requires comparing the model output with reference translations. Compared with lexical overlap based methods such as BLEU, neural metrics like BLEURT and COMET can capture semantic relationships. Reference-free evaluation (OpenKiwi (Kepler et al., 2019), TransQuest (Ranasinghe et al., 2020), COMET-QE (Rei et al., 2022b) and MetricX-QE (Juraska et al., 2023)), also named quality estimation (QE), mainly indicates that directly predicting a quality score for the translation based on the input and model output.

Adding tags is an effective method to integrate additional signals to neural models, such as formality level (Yamagishi et al., 2016), politeness (Sennrich et al., 2016a), domain (Kobus et al., 2017), and so on. Johnson et al. (2017) indicate the target language for multilingual NMT. Scarton and Specia (2018) indicate the audience for text simplification. Bandel et al. (2022) control paraphrasing via semantic similarity, syntactic and lexical variation.

5 Conclusion

In this paper, we facilitate quality-aware NMT by self-estimated quality with bi-directional translation losses and representing quality scores via vector interpolation. Experiments show that our method can lead to consistent and significant improvements on low/middle/high-resource tasks. Our analysis shows that the model is aware of the quality after training and can generate translations of the given quality score.

Acknowledgements

We appreciate our reviewers for their insightful comments and suggestions. This work is partially supported by the National Natural Science Foundation of China (Grant No. 62306284), China Postdoctoral

Science Foundation (Grant No. 2023M743189), and the Natural Science Foundation of Henan Province (Grant No. 232300421386).

References

- Ali Araabi and Christof Monz. 2020. Optimizing transformer for low-resource neural machine translation. In Donia Scott, Nuria Bel, and Chengqing Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3429–3435, Barcelona, Spain (Online), December. International Committee on Computational Linguistics.
- Elron Bandel, Ranit Aharonov, Michal Shmueli-Scheuer, Ilya Shnayderman, Noam Slonim, and Liat Ein-Dor. 2022. Quality controlled paraphrase generation. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 596–609, Dublin, Ireland, May. Association for Computational Linguistics.
- Ondrej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, and Aleš Tamchyna. 2014. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 12–58, Baltimore, Maryland, USA, June. Association for Computational Linguistics.
- Mauro Cettolo, Jan Niehues, Sebastian Stüker, Luisa Bentivogli, and Marcello Federico. 2014. Report on the 11th IWSLT evaluation campaign. In Marcello Federico, Sebastian Stüker, and François Yvon, editors, *Proceedings of the 11th International Workshop on Spoken Language Translation: Evaluation Campaign*, pages 2–17, Lake Tahoe, California, December 4–5.
- Vishrav Chaudhary, Yuqing Tang, Francisco Guzmán, Holger Schwenk, and Philipp Koehn. 2019. Low-resource corpus filtering using multilingual sentence embeddings. In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2)*, pages 261–266, Florence, Italy, August. Association for Computational Linguistics.
- Mia Xu Chen, Orhan Firat, Ankur Bapna, Melvin Johnson, Wolfgang Macherey, George Foster, Llion Jones, Mike Schuster, Noam Shazeer, Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Zhifeng Chen, Yonghui Wu, and Macduff Hughes. 2018. The best of both worlds: Combining recent advances in neural machine translation. In Iryna Gurevych and Yusuke Miyao, editors, *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 76–86, Melbourne, Australia, July. Association for Computational Linguistics.
- Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States, July. Association for Computational Linguistics.
- Marina Fomicheva, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. 2020. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:539–555.
- Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.
- Marcin Junczys-Dowmunt. 2018. Microsoft’s submission to the WMT2018 news translation task: How I learned to stop worrying and love the data. In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 425–430, Belgium, Brussels, October. Association for Computational Linguistics.
- Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In Philipp Koehn, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 756–767, Singapore, December. Association for Computational Linguistics.

- Fabio Kepler, Jonay Trénous, Marcos Treviso, Miguel Vera, and André F. T. Martins. 2019. OpenKiwi: An open source framework for quality estimation. In Marta R. Costa-jussà and Enrique Alfonseca, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 117–122, Florence, Italy, July. Association for Computational Linguistics.
- Catherine Kobus, Josep Crego, and Jean Senellart. 2017. Domain control for neural machine translation. In Ruslan Mitkov and Galia Angelova, editors, *Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017*, pages 372–378, Varna, Bulgaria, September. INCOMA Ltd.
- Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In Philipp Koehn, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- T. KUDO. 2010. Mecab : Yet another part-of-speech and morphological analyzer. <http://mecab.sourceforge.net/>.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Pierre Isabelle, Eugene Charniak, and Dekang Lin, editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July. Association for Computational Linguistics.
- Jan-Thorsten Peter, David Vilar, Daniel Deutsch, Mara Finkelstein, Juraj Juraska, and Markus Freitag. 2023. There’s no data like better data: Using QE metrics for MT data filtering. In Philipp Koehn, Barry Haddow, Tom Kocmi, and Christof Monz, editors, *Proceedings of the Eighth Conference on Machine Translation*, pages 561–577, Singapore, December. Association for Computational Linguistics.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium, October. Association for Computational Linguistics.
- Tharindu Ranasinghe, Constantin Orasan, and Ruslan Mitkov. 2020. TransQuest: Translation quality estimation with cross-lingual transformers. In Donia Scott, Nuria Bel, and Chengqing Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5070–5081, Barcelona, Spain (Online), December. International Committee on Computational Linguistics.
- 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. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In Philipp Koehn, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In Philipp Koehn, Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri, editors, *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid), December. Association for Computational Linguistics.
- Carolina Scarton and Lucia Specia. 2018. Learning simplifications for specific target audiences. In Iryna Gurevych and Yusuke Miyao, editors, *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 712–718, Melbourne, Australia, July. Association for Computational Linguistics.

- Julian Schamper, Jan Rosendahl, Parnia Bahar, Yunsu Kim, Arne Nix, and Hermann Ney. 2018. The RWTH Aachen University supervised machine translation systems for WMT 2018. In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo Negri, Aurélie Névél, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Verspoor, editors, *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 496–503, Belgium, Brussels, October. Association for Computational Linguistics.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online, July. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016a. Controlling politeness in neural machine translation via side constraints. In Kevin Knight, Ani Nenkova, and Owen Rambow, editors, *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 35–40, San Diego, California, June. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016b. Neural machine translation of rare words with subword units. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August. Association for Computational Linguistics.
- Qwen Team. 2025. Qwen3 technical report.
- Christian Tomani, David Vilar, Markus Freitag, Colin Cherry, Subhajit Naskar, Mara Finkelstein, Xavier Garcia, and Daniel Cremers. 2024. Quality-aware translation models: Efficient generation and quality estimation in a single model. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15660–15679, Bangkok, Thailand, August. Association for Computational Linguistics.
- 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 I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.
- Hongfei Xu, Qiuhui Liu, Josef van Genabith, Deyi Xiong, and Meng Zhang. 2021. Multi-head highly parallelized LSTM decoder for neural machine translation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 273–282, Online, August. Association for Computational Linguistics.
- Hayahide Yamagishi, Shin Kanouchi, Takayuki Sato, and Mamoru Komachi. 2016. Controlling the voice of a sentence in Japanese-to-English neural machine translation. In Toshiaki Nakazawa, Hideya Mino, Chenchen Ding, Isao Goto, Graham Neubig, Sadao Kurohashi, Ir. Hammam Riza, and Pushpak Bhattacharyya, editors, *Proceedings of the 3rd Workshop on Asian Translation (WAT2016)*, pages 203–210, Osaka, Japan, December. The COLING 2016 Organizing Committee.