

Challenging the state-of-the-art Machine Translation Metrics from a Linguistic Perspective

Eleftherios Avramidis, Shushana Manakhimova, Vivien Macketanz and Sebastian Möller

German Research Center for Artificial Intelligence (DFKI),

Speech and Language Technology, Berlin, Germany

firstname.lastname@dfki.de

Abstract

We employ a linguistically motivated challenge set in order to evaluate the state-of-the-art machine translation metrics submitted to the Metrics Shared Task of the 8th Conference for Machine Translation. The challenge set includes about 21,000 items extracted from 155 machine translation systems for three language directions (German \leftrightarrow English, English \rightarrow Russian), covering more than 100 linguistically-motivated phenomena organized in 14 categories. The metrics that have the best performance with regard to our linguistically motivated analysis are the COMETOID22-WMT23 (a trained metric based on distillation) for German-English and METRICX-23-C (based on a fine-tuned mT5 encoder-decoder language model) for English-German and English-Russian. Some of the most difficult phenomena are *passive voice* for German-English, *named entities*, *terminology* and *measurement units* for English-German and *focus particles*, *adverbial clause* and *stripping* for English-Russian.

1 Introduction

Most NLP evaluation has relied for years on testing the system performance on randomly picked test sets and producing a single generic score. Yet, machine learned systems learn to make abstractions and due to these, phenomena who are on the long tail of the training and test data may be overlooked hidden behind a very high generic score. Additionally, generic scores are often helpful to show relative improvement and reflect overall quality, but cannot explain the performance in a comprehensive way.

For example, old-style machine translation (MT) metrics measuring lexical overlap would equally penalize the omission of an article and the omission of the particle forming the negation in a sentence, although negation is more crucial for its meaning. While the evaluation of so obvious errors has been

addressed by the trained MT metrics, their evaluation relies on correlations with human judgments on randomly picked test-sets. In this case, a single correlation score may not be able to explain the strengths and weaknesses of the metrics with regard to the functioning of language.

Motivated by these considerations, we employ a multifold test set with linguistically-motivated challenges that will allow us to understand the metric performance from a linguistic perspective. These challenges are organized in smaller sets, one set per phenomenon, whereas the phenomena are organized in broader categories. By measuring the ability of the metrics to detect the errors in these challenge sets, we can get scores that indicate different aspects of linguistic performance.

This paper describes the application of such a challenge set on the evaluation of the MT metrics submitted at the relevant shared task of the 8th Conference of Machine Translation ([Freitag et al., 2023](#)). The rest of the paper is structured as following: Section 2 describes related work, and section 3 describes the way the challenges were selected. In Section 4 the results are presented and described, first from the perspective of metric comparison and then focusing on the performance for particular linguistically-motivated categories and phenomena per language direction. Some conclusions are given in Section 5.

2 Related work

There has been a growing interest for more fine-grained evaluation of Natural Language Processing (NLP) tools, as shown by the increasing number of publications many of whom have received distinctions ([Ribeiro et al., 2020](#); [Avelino et al., 2022](#); [Campolungo et al., 2022](#)). Concerning machine translation (MT), initial efforts were made in the 1990s with the introduction of test suites ([King and Falkedal, 1990](#)), and these efforts have been revitalized in light of recent advancements in the

field (Guillou and Hardmeier, 2016). To the best of our knowledge, the first endeavours related to the use of challenge sets in a meta-level in order to evaluate MT metrics were applied to Quality Estimation metrics (Avramidis et al., 2018), based on the first version of our linguistically-motivated test suite (Macketanz et al., 2018). The analysis was broadened to cover a broader range of MT metrics, including reference-based ones, as appeared in the Findings paper of the Metrics shared task of the 6th Conference on Machine Translation (Freitag et al., 2021), which was based on a later version of our test suite on German-English (Avramidis et al., 2019, 2020; Macketanz et al., 2021, 2022a), a resource also employed in this paper.

With the occasion of the first challenge set sub-task for the metrics shared task of the 7th Conference on Machine Translation (Freitag et al., 2022), a few more challenge sets emerged. ACES (Amrhein et al., 2022) for example, focuses on 68 accuracy errors. Similarly, Alves et al. (2022) evaluate the robustness of MT metrics by generating translations with critical errors. In a more linguistic direction, Chen et al. (2022) examine the capability of the metrics to correlate synonyms in different areas and to discern catastrophic errors at both word- and sentence-levels.

Our submission at that sub-task (Avramidis and Macketanz, 2022) augmented the preliminary analysis appearing at Freitag et al. (2021) by adding the language direction of English-German and presenting a more fine-grained analysis, not only in the category level but also on the phenomenon level. This year’s submission, explained on our paper, includes that same challenge set as last year, whereas English-Russian has been added as an additional language direction.

3 Method

3.1 Test suite for MT systems

Here, we are going to explain how we created the pool of MT sentences that were used for the challenge set. The selection was based on a linguistically-motivated test suite (Macketanz et al., 2022a)¹. The test suite contains a set of source sentences focusing on particular phenomena, each of them accompanied by some rules or regular expressions that can detect which translations would be accepted for these source sentences. This allows a

semi-automatic evaluation when new translations are provided, whereas a human annotator resolves cases not covered by the rules.

For this experiment, we employed the test suite on three language directions: German-English (Avramidis et al., 2020), English-German (Macketanz et al., 2021) and English-Russian (Macketanz et al., 2022b). The German-English side consists of 5,539 German test sentences covering 107 linguistically motivated phenomena, the English-German side consists of 4,782 English test sentences covering 126 phenomena, and the English-Russian side consists of 1,225 English test sentences covering 64 phenomena. All language directions are organized in 14 categories, which nevertheless differ among the directions.

The above described test suite has been used to evaluate the outputs of 116 German-English, 29 English-German systems and 10 English-Russian systems submitted at the translation task of the Conference of Machine Translation (WMT). German-English outputs were collected from systems submitted in the years 2018-2021, English-German outputs in the years 2020-2021 and English-Russian in 2022.

3.2 Challenge set for MT metrics

The sentences selected with the help of the test suite are consequently used to create the challenge set. The source sentences and the system outputs have to be organized in contrastive pairs of correct/incorrect translations and a reference. In order to achieve this, for every source sentence from the test suite selection we create a challenge item including:

- one correct translation to be used as a reference translation,
- another correct translation to be used as the first translation candidate
- one incorrect translation to be used as the contrastive translation candidate

The two candidate translations and the reference consist one challenge item. Since source and translations were collected as a result of testing for a particular phenomenon, the same phenomenon will be what the challenge item will test.

Given that we may have many correct and wrong translations for the same source, the reference and the translations of the challenge items result from random combinations of correct and wrong translations from the collected WMT outputs. Therefore,

¹<https://github.com/DFKI-NLP/mt-testsuite>

the same source sentence may appear many times.

As a result, we get a challenge set with 10,402 items for German-English, 8,945 items for English-German and 1,727 items for English-Russian.

3.3 Evaluation of metrics

For each challenge item, the two machine translation (MT) outputs, are provided to the metrics as separate MT hypotheses. Which output is correct, and which is incorrect, is hidden from the metrics. These hypotheses are then evaluated against the previously mentioned reference and/or the source. An item is deemed correctly scored when the metric assigns a higher score to the correct MT output compared to the incorrect one. Following this, the statistics below are computed:

- i) **Accuracy per Phenomenon:** the ratio of all correctly-scored challenge items per phenomenon to the total number of challenge items for that particular phenomenon.
- ii) **Accuracy per Category:** the ratio of all correctly-scored challenge items per category to the total number of challenge items for that category, after consolidating the underlying phenomena of that category into a single set.

Significance tests are performed to compare the highest metric accuracy for each phenomenon with all other metric accuracies for the same phenomenon. This is a one-tailed Z-test, conducted with a significance level of $\alpha = 0.95$. Metrics with accuracies that are not significantly worse than the highest accuracy are considered to share the top position for that phenomenon. A similar approach is used to identify the best accuracies per category, after aggregating the challenge items from the underlying phenomena within each category.

Metric categories We conduct this significance testing in two stages: first, for all metrics involved in the shared task, and then separately for each of the three metric categories (baseline, Quality Estimation (QE) as a metric, reference-based). Systems that are significantly superior per phenomenon across all metrics are highlighted with a gray background, while those that are significantly superior per metric category are denoted in boldface.

Averaging Lastly, we provide three types of averaging scores:

- i) **Micro-average:** This approach considers all items equally, aggregating all test items to compute the average percentages.
- ii) **Category macro-average:** Here, all categories are treated equally, with the percentages being computed independently for each category and then averaged.
- iii) **Phenomenon macro-average:** This average treats all phenomena equally, with the percentages being computed independently for each phenomenon and then averaged.

4 Results

The results are displayed in detail in Tables 1, 2 and 3 for the category level and in Tables 4, 5 and 6 for the phenomenon level, for the three language directions respectively.

4.1 Metric performance analysis

Here we are observing the statistics with a focus on comparing the performance of various metrics on the challenge set.

German-English The accuracies of the metrics, as measured for several categories in German-English, can be seen in Table 1. The best performing metric for German-English is COMETOID22-WMT23 (Gowda et al., 2023), which, wins significantly based on both the micro-average (83%) and the macro-average (87%). This metric is a distilled QE model that has been trained on COMET (Rei et al., 2020) scores of WMT outputs, including the ones of WMT23. For this reason, we include it into the reference-aware metrics. We notice that its performance among the other metrics is impressive. It is the first metric in 6 categories and among them the only one who wins at *Verb tense/aspect/mood* and *function words*, achieving 93% and 91% accuracy respectively.

Another two reference-based baseline metrics, COMET and PRISMREF (Thompson and Post, 2020a,b) share the first position when the category macro-average is considered (82%). None of the other reference-aware metrics submitted this year managed to compete with the metrics with the highest accuracy mentioned above.

The lowest performing metric is the referenceless random baseline RANDOM-SYSNAME, provided by the organizers (44%), followed by XLSIMQE (55-58%; Mukherjee and Shrivastava,

2023) and MATESE (57-58%; Perrella et al., 2022).

When considering the metric performance with regard to particular categories, one can see, again this year, that different metrics win in different combinations of categories. Here, only COMETOID22-WMT23 as mentioned above, wins 6 metrics, followed by PRISMREF and METRIC-23-C, which win 4 categories. 17 metrics do not win any category, ranging in accuracies around 75%.

English-German The accuracies of the metrics, as measured for several categories in English-German, can be seen in Table 2. The best performing metric in English-German is METRICX-23-C, which is in the first significance cluster based on both the micro-average (81%) and the category macro-average (84%). This metric uses the mT5 encoder-decoder language model, which is fine-tuned using direct assessment data, MQM (Lommel et al., 2014) data and synthetic data. The categories to which its success may be mostly attributed are the *multi-word expressions* (MWE; with 92%) and the non-verbal agreement (95%).

Another three metrics share the first position, when the micro average is considered, namely the QE version of the latter, *MetricX-23-QE-c* and also *mbr-metricx-qe* (Naskar et al., 2023) and XCOMET-Ensemble. It is impressive that QE methods manage to reach high accuracy without access to reference content.

When looking at the worse-performing metrics, MATESE here performs worse than the baseline (36-38%), followed by PARTOKENGRAM_F (55-56%; Dreano et al., 2023b).

In English-German it is even harder to say which metrics perform well in multiple categories, as only one of them, XCOMET-QE-ENSEMBLE, achieves the highest performance in 3 categories (*function words, non-verbal agreement and subordination*). The rest of the metrics show a good performance in 2 categories or fewer.

English-Russian The accuracies of the metrics, as measured for several categories in English-German, can be seen in Table 3. For this language pair, variants of the MetricX achieve significantly higher accuracies than all the other metrics. In particular, METRICX-23-C achieves the highest accuracy based on both micro-average and category macro-average, whereas METRICX-23-B and METRICX-23-QE-C achieve a slightly

lower macro-average, which is nevertheless not significantly worse than the one of the former. MATESE is again by far the lowest performing metric (32/34%), lower than the random baseline. We may assume that this metric has not been optimized for this language direction.

4.2 Linguistically motivated analysis

In this section, we are focusing on the results for particular phenomena or categories.

4.2.1 German-English

Category-level The overall average accuracy of all metrics with regard to the linguistically motivated categories is at 76% for German-English, which is two percentage points lower than last year's average. It is still a fact, that the metrics in average fail to predict properly the scores for one out of four challenge items that we provided. Luckily, there has been noticeable accuracy for some categories, for example METRICX reference-based variants achieved an accuracy of 96% for *false friends*, whereas *negation* errors have been scored correctly with a 98.5

The worse performing category is *Verb valency*, where the best metrics achieved only 66% accuracy, and the rest of the metrics averaged to a mere 56%. In this category one can observe the lowest accuracy, given by an LLM-based metric, EM-BED_LLAMA (Dreano et al., 2023a) with 41%.

Phenomenon-level The best accuracy for this language pair (Table 4) is achieved this year at several variations of verb tenses, i.e. *Transitive - future II, Modal negated - present, Reflexive - preterite subjunctive II* and *Intransitive - pluperfect* which get more than 85% in average.

The lowest accuracy of all metrics in average is given for *passive voice*, where the highest accuracy achieved by several metrics is only 54.5%. Errors related to *commas, domain-specific terms* and *locations* have also been scored with a less than 65% accuracy.

4.2.2 English-German

Category-level The overall average accuracy of all metrics with regard to the linguistically motivated categories is at 71-73% for English-German. The category where all metrics perform better in average is *negation* (83%), where 11 of the metrics achieve more than 90% accuracy. Negation is closely followed by *function words Non-verbal agreement* (80%).

The worse performing category in average is *named entities and terminology* (58,8%), where most metrics' accuracies are close to 50%, except for BLEURT (Yan et al., 2023) (80.3%). The rest of the categories lie in rather mediocre accuracies, between 58.8% and 80%.

Phenomenon-level The English-German phenomena, where metrics perform best in average (Table 5) are the *transitive conditional II simple, gerunds, contact clause* and the *intransitive present perfect simple*, achieving more than 85% of accuracy. The phenomena which incur the lowest average accuracies are the *transitive present progressive, measuring units, modals* and *intransitive -future II progressive* with less than 50% accuracy. The former and the latter were observed as the most difficult phenomena to score also last year.

4.3 English-Russian

This analysis for English-Russian occurs for the first time this year, based on the MT outputs collected at last year's shared task. For this purpose the test instances are much fewer than the other language pairs and therefore the numbers are not very conclusive. Therefore, categories and phenomena that have only a handful of samples will not be included in our analysis, although they appear in the tables.

Category-level Here, the average accuracy over all metrics is much lower than the other language directions, reaching only 66%, only 20% above the random baseline. The best performing category is *ambiguity* (86,3%), more than 13% better than the following categories. The worst performing categories are *function words* and *punctuation*, with less than 55%. The rest of the categories range in accuracies between 53 and 73%.

Phenomenon-level The good performance of the *ambiguity* category is also confirmed in the table on the phenomenon level (Table 6), as in Russian this is the only phenomenon of this category, as opposed to other language pairs where we also have examples of structural ambiguity. The most difficult phenomena to score appear to be the *focus particles, adverbial clause* and *stripping* with less than 50% average accuracy, in many cases lower than the random baseline.

5 Conclusion

In this paper we analysed the performance of several state-of-the-art metrics with regard to particular linguistically-motivated phenomena for three language pairs, German-English, English-German and for the first time, English Russian. The analysis gave a multitude of observations, regarding both the performance of the metrics and the corresponding linguistic observations.

The metrics demonstrating the best performance in average are COMETOID22-WMT23 for the German-English language pair, and METRICX-23-C for both the English-German and English-Russian language pairs. Quality estimation methods have impressively good performance in several phenomena. Some metrics that report usage of LLMs (EMBED_LLAMA) have not scored very high in overall, indicating that more work is required in this direction.

Among the various linguistic phenomena, we could identify some of the particularly challenging ones. In German-English, metrics have difficulties scoring the *passive voice* properly. In English-German *named entities and terminology* as well as specific *measurement units* pose the most difficulties. In English-Russian translation, translations with *focus particles, adverbial clause*, and *stripping* phenomena emerge as particularly complex challenges.

Acknowledgements

This research was supported by the German Research Foundation (Deutsche Forschungsgemeinschaft; DFG) through the project TextQ and by the German Federal Ministry of Education (BMBF) through the project SocialWear (grant num. d01IW20002). We would like to thank Hans Uszkoreit, Aljoscha Burchardt, Ursula Strohriegel, Renlong Ai, He Wang, Ekaterina Lapshinova-Koltunski and Sergei Bagdasarov for their prior contributions for the creation of the test suite.

References

- Duarte Alves, Ricardo Rei, Ana C Farinha, José G. C. de Souza, and André F. T. Martins. 2022. Robust MT evaluation with sentence-level multilingual augmentation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 469–478, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

- Chantal Amrhein, Nikita Moghe, and Liane Guillou. 2022. **ACES: Translation accuracy challenge sets for evaluating machine translation metrics**. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 479–513, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Mariana Avelino, Vivien Macketanz, Eleftherios Avramidis, and Sebastian Möller. 2022. **A Test Suite for the Evaluation of Portuguese-English Machine Translation**. In *Computational Processing of the Portuguese Language*, pages 15–25, Cham. Springer International Publishing.
- Eleftherios Avramidis and Vivien Macketanz. 2022. **Linguistically motivated evaluation of machine translation metrics based on a challenge set**. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 514–529, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Eleftherios Avramidis, Vivien Macketanz, Arle Lommel, and Hans Uszkoreit. 2018. **Fine-grained evaluation of quality estimation for machine translation based on a linguistically motivated test suite**. In *Proceedings of the AMTA 2018 Workshop on Translation Quality Estimation and Automatic Post-Editing*, pages 243–248, Boston, MA. Association for Machine Translation in the Americas.
- Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, Aljoscha Burchardt, and Sebastian Möller. 2020. **Fine-grained linguistic evaluation for state-of-the-art machine translation**. In *Proceedings of the Fifth Conference on Machine Translation*, pages 346–356, Online. Association for Computational Linguistics.
- Eleftherios Avramidis, Vivien Macketanz, Ursula Strohriegel, and Hans Uszkoreit. 2019. **Linguistic evaluation of German-English machine translation using a test suite**. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 445–454, Florence, Italy. Association for Computational Linguistics.
- Niccolò Campolungo, Federico Martelli, Francesco Saina, and Roberto Navigli. 2022. **DiBiMT: A novel benchmark for measuring Word Sense Disambiguation biases in Machine Translation**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4331–4352, Dublin, Ireland. Association for Computational Linguistics.
- Xiaoyu Chen, Daimeng Wei, Hengchao Shang, Zongyao Li, Zhanglin Wu, Zhengze Yu, Ting Zhu, Mengli Zhu, Ning Xie, Lizhi Lei, Shimin Tao, Hao Yang, and Ying Qin. 2022. **Exploring robustness of machine translation metrics: A study of twenty-two automatic metrics in the WMT22 metric task**. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 530–540, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Sören Dreano, Derek Molloy, and Noel Murphy. 2023a. **Embed_Llama: using LLM embeddings for the Metrics Shared Task**. In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, Singapore, Singapore (Hybrid). Association for Computational Linguistics.
- Sören Dreano, Derek Molloy, and Noel Murphy. 2023b. **Tokengram_F, a fast and accurate token-based chrF++ derivative**. In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, Singapore, Singapore (Hybrid). Association for Computational Linguistics.
- Markus Freitag, Nitika Mathur, Chi-ku Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thomson, Tom Kocmi, Frédéric 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 (WMT)*, Singapore, Singapore. Association for Computational Linguistics.
- Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-ku Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. **Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust**. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-ku Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021. **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*, pages 733–774, Online. Association for Computational Linguistics.
- Thamme Gowda, Tom Kocmi, and Marcin Junczys-Dowmunt. 2023. **Cometoid: Distilling Strong Reference-based Machine Translation Metrics into Even Stronger Quality Estimation Metrics**. In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, Singapore, Singapore (Hybrid). Association for Computational Linguistics.
- Liane Guillou and Christian Hardmeier. 2016. **PROTEST: A Test Suite for Evaluating Pronouns in Machine Translation**. In *Tenth International Conference on Language Resources and Evaluation (LREC 2016)*.
- Margaret King and Kirsten Falkedal. 1990. **Using test suites in evaluation of machine translation systems**.

- In *Proceedings of the 13th conference on Computational Linguistics*, volume 2, pages 211–216, Morristown, NJ, USA. Association for Computational Linguistics.
- Arle Lommel, Aljoscha Burchardt, Maja Popović, Kim Harris, Eleftherios Avramidis, and Hans Uszkoreit. 2014. [Using a new analytic measure for the annotation and analysis of MT errors on real data](#). In *Proceedings of the 17th Annual Conference of the European Association for Machine Translation*, pages 165–172, Dubrovnik, Croatia. European Association for Machine Translation.
- Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, and Hans Uszkoreit. 2018. [Fine-grained evaluation of German-English machine translation based on a test suite](#). In *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*, pages 578–587, Belgium, Brussels. Association for Computational Linguistics.
- Vivien Macketanz, Eleftherios Avramidis, Aljoscha Burchardt, He Wang, Renlong Ai, Shushen Manakhimova, Ursula Strohriegel, Sebastian Möller, and Hans Uszkoreit. 2022a. [A linguistically motivated test suite to semi-automatically evaluate German-English machine translation output](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 936–947, Marseille, France. European Language Resources Association.
- Vivien Macketanz, Eleftherios Avramidis, Shushen Manakhimova, and Sebastian Möller. 2021. [Linguistic evaluation for the 2021 state-of-the-art machine translation systems for German to English and English to German](#). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1059–1073, Online. Association for Computational Linguistics.
- Vivien Macketanz, Shushen Manakhimova, Eleftherios Avramidis, Ekaterina Lapshinova-koltunski, Sergei Bagdasarov, and Sebastian Möller. 2022b. [Linguistically motivated evaluation of the 2022 state-of-the-art machine translation systems for three language directions](#). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 432–449, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Ananya Mukherjee and Manish Shrivastava. 2023. [MEE4 and XLSim : IIIT HYD's Submissions' for WMT23 Metrics Shared Task](#). In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, Singapore, Singapore (Hybrid). Association for Computational Linguistics.
- Subhajit Naskar, Daniel Deutsch, and Markus Freitag. 2023. [Quality Estimation using Minimum Bayes Risk](#). In *Proceedings of the Eighth Conference on Machine Translation (WMT)*, Singapore, Singapore (Hybrid). Association for Computational Linguistics.
- Stefano Perrella, Lorenzo Proietti, Alessandro Scirè, Niccolò Campolungo, and Roberto Navigli. 2022. [MaTESe: Machine translation evaluation as a sequence tagging problem](#). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 569–577, Abu Dhabi, United Arab Emirates (Hybrid). 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.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. [Beyond accuracy: Behavioral testing of NLP models with CheckList](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.
- Brian Thompson and Matt Post. 2020a. [Automatic machine translation evaluation in many languages via zero-shot paraphrasing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 90–121, Online. Association for Computational Linguistics.
- Brian Thompson and Matt Post. 2020b. [Paraphrase generation as zero-shot multilingual translation: Disentangling semantic similarity from lexical and syntactic diversity](#). In *Proceedings of the Fifth Conference on Machine Translation*, pages 561–570, Online. Association for Computational Linguistics.
- Yiming Yan, Tao Wang, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Mingxuan Wang. 2023. [BLEURT has universal translations: An analysis of automatic metrics by minimum risk training](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5428–5443, Toronto, Canada. Association for Computational Linguistics.

Table 1: Accuracy of the metrics (%) with regard to the 14 linguistically motivated categories for German-English. The significantly best systems per phenomenon over all metrics are indicated with a gray background, whereas the significantly best systems per metrics category are indicated with boldface.

ling. category	#	BERT score	BLEURT-20	COMET	YISI-1	chRF	spBLEU	COMET-QE	Calibri-COMET22-QE	CometKiwix-XL	CometKiwix-XXL	GEMBA-MQM	KG-BERT	MS-COMET-QE	MetricX-23-QE-b	MetricX-23-QE-c	MetricX-23-QE-e	XLSimQE	cometoid22-wmt21	cometoid22-wmt22	cometoid22-wmt23	mbr-metricx-qe0p2pl-qe	prismSrc2	Calibri-COMET22	MEE4_stsb_xlm	MATTS-e	MetricX-23-b	MetricX-23-c	XCOMET-XXL	XCOMET-XL	XLsim	eBLU	embed_Llama	mDlregressor	partkengram_F	tokengram_F	avg		
baseline		QE as a metric												ref. based metrics																									
Ambiguity	146	84 71 90 84 86 88	97 88	42 77 60 42 77	42 39	75	85	82 46	56 38	59	58	87	86	95 14	84 90	92 50	88	95 88	82	88	62	93 70	62	81	62	88	73												
Coordination & ellipsis	836	69 61	80 78 69 61	74 63	73	74 75	74 60	74 71	78	79	80 48	77 38	75	74	76	83 80	58 75	62	63	33 76	81 72	79	81 64	75	61	56	72	49	62	69									
False friends	225	65 63	69 74 71 71	67 65	71	61 75	71 65	71 73	64	80	66 46	76 82	64	52	66	76 63	68	74	83	87 30	60	82	62	65	59	53	76	66	55	56	52	68	67						
Function word	200	90	78 79	90 82	74	93 74	92	94 92	90	82	90	78	86	72	85	36	95 62	94 94	90	67	84	86	82	76	60	86	84	94 86	64	76	70	52	76	60	76	80			
MWE	829	79 72	87 90	86 77	85 74	74	86	79	73	66	73	71	88	89	85	45	83	37	78	90 88	89	76	78	18	86	92 85	90	93 76	81	69	69	72	54	78	76				
Named entity & termin.	1272	57	54	66 64	62	61	67 62	54	51	49	55	61	65	58	44	54	46	60	54	68	62	80 56	62	60	63	30	72	74 71	62	60	60	71	60	48	54	49	60	59	
Negation	174	87	83	89 92	84	84	86	90	90	84	89	76	89	92	83	88	78	47	84	40	90	90	96 89	91	91 58	85	91 80	84	85	79	80	86	67	88	52	86	83		
Non-verbal agreement	372	74 72	81	88 78	70	83 75	82	90	84	81	83	81	78	95 93	91	45	95 43	93	91	93	82	94	71	87	76	69	57	90	95 87	94 91	88	73	72	65	84	54	69	80	
Punctuation	336	69	73	74	70	66	72	77 69	82 72	76	82 60	82 65	73	66	73	42	60	44	81	82 81	75	70	75	65 74	72	28	74 70	67	65	65	46	68	74 49	57	44	73	68		
Subordination	994	78	74	81	84	78	75	86 74	89 88	89 76	89 80	83	85	83	42	89 43	83	83	86	83	84	76	83	79	78	16	83	84	83	89 85	79	75	75	71	75	56	76	78	
Verb tense/aspect/mood	3081	68	62	70 70	69	67	64	71	72	75	75	71	75	61	82 84	85 43	83	52	62	65	75	71	77	61	70	74	46	77	78	83 76	71	69	72	55	71	62	69	70	74
Verb valency	480	73	64	82 79	74	70	79	69	77	81	85	77	62	77	71	86	80	86	42	87	52	79	80	77	87 66	77	70	70	36	84 81	85 80	75	79	66	69	65	70	70	74
macro avg.	8945	74	69	79 80	76	73	80 72	75	78	77	75	70	75	75 82	81 81	64	79	76	76	38	80	84 79	81	79	68	76	70	60	71	55	73	73							
micro avg.	8945	70	65	75 76	72	69	74	68	73	74	75	74	68	74	67	79	81 80	44	78	47	71	71	78	77	81 62	74	71	73	36	79	81 78	81 77	70	73	69	56	70	71	

Table 2: Accuracy of the metrics (%) with regard to the 12 linguistically motivated categories for English-German

		ling category		baseline		QE as a metric		ref. based metrics		avg	
#	BERTscore	BLEU	BLEURT-20	COMET	YISI-1	chrF	spBLEU	COMET	Cometrikiwi-XL	Cometrikiwi-XXL	Cometrikiwi
66	80 67	100 97	88 67	83 73	100	94	100	100	95	100	100
203	84 73	68 73	76 67	77 69	73	64	68	38	68 76	62 79	64 45
6	100 17	100 100	100	17 83	50	0	0	0	0 100	0 17	17 17
116	59 69	69 59	53 38	48 72	53	55	50	51	50 36	66 55	48 44
122	70 64	79 80	74 64	79 67	72	70	77	57	77 77	75 88	75 40
243	79 69	95 91	79 76	85 76	84	62	59	54	59 69	63 73	73 49
34	74 74	79 68	74 85	68 74	68	85	91	85	41 71	100 95	91 9
61	69 56	82 77	67 67	61 57	84	100	54	72	54 74	72 98	66 41
121	42 46	53 68	86 41	73 55	66	29	65	27	65 87	55 66	43 41
499	62 62	61 70	67 66	63 63	72	66	41	66 75	65 66	61 45	77 32
135	79 62	76 85	76 84	81 74	79	73	61	73 64	67 65	66 39	74 60
121	70 67	70 75	76 67	56 68	71	88	79	66	79 71	76 85	76 36
1727	72 60	78 79	76 61	71 66	69	66	65	63	65 64	64 81	58 41
1727	69 64	72 76	73 65	70 67	74	67	68	51	68 71	67 74	64 44
macro avg.		77 44	59 70	88 77	72 62	65 64	66 74	42 34	81 83	79 74	62 63
micro avg.		77 55	71 71	83 82	73 77	68 74	73 71	40 32	73 73	71 74	55 57

Table 3: Accuracy of the metrics (%) with regard to the linguistically motivated categories for English-Russian

Accuracy of the metrics (%) with regards to the linguistically-motivated phenomena for German-English - Continued on next page

Accuracy of the metrics (%) with regards to the linguistically-motivated phenomena for German-English - Continued on next page

Table 4: Accuracy of the metrics(%) with regard to the linguistically v-motivated phenomena for German-English

	ling. phenomenon	BERTScore	BLEU	BLURT-20	COMET	VIS-I	chF	spBERT	QE as a metric	ref. based metrics	avg
#											
Ambiguity											
Lexical ambiguity	146	84	71	90	84	86	88	97	88	42	77
Coo ordination & ellipsis											
Gapping	163	72	59	86	84	71	67	88	66	75	64
Pseudogapping	201	81	76	98	92	79	67	81	73	94	97
Right node raising	47	83	64	87	98	83	72	89	70	91	85
Slicing	169	55	58	57	50	54	45	58	47	53	39
Stripping	139	72	58	68	76	68	58	76	61	64	36
VP-ellipsis	117	59	47	85	84	68	53	74	45	88	85
False friends											
False friends	225	65	63	69	74	71	67	65	71	64	60
Function word											
Focus particle	20	45	30	70	60	50	35	70	30	40	40
Question tag	180	96	83	80	93	86	78	96	79	90	77
MWE											
Collocation	112	73	62	93	90	89	77	88	60	96	92
Compound	63	73	51	56	89	95	84	98	68	97	100
Idiom	266	86	82	97	98	93	75	83	93	90	90
Nominal MWE	288	81	72	83	86	71	78	75	74	78	75
Propositional MWE	35	69	86	71	83	83	83	83	80	83	81
Verbal MWE	65	68	72	88	77	66	74	78	69	75	66
Named entity & termin.											
Date	234	52	55	74	71	64	60	68	63	67	77
Domain-specific term	312	69	56	88	82	79	74	88	66	71	67
Location	12	67	83	75	67	50	58	100	83	67	72
Measuring unit	389	54	48	53	56	55	56	68	53	37	28
Proper name	325	50	57	54	50	53	54	56	66	50	44
Negation											
Negation	174	87	83	89	92	92	84	86	90	89	87
Non-verbal agreement											
Correlation	81	85	85	95	80	73	86	84	88	83	77
Genitive	206	75	73	93	82	68	80	76	77	94	73
Possession	85	62	55	86	86	74	58	87	60	93	93
Punctuation											
Quotation marks	336	69	73	74	70	66	72	77	69	82	76
Subordination											
Adverbial clause	193	72	80	82	71	73	83	77	91	90	91
Cleft sentence	179	67	64	61	75	66	59	75	62	74	75
Contact clause	150	84	75	95	91	75	89	75	97	93	91
Indirect speech	38	55	42	63	47	82	42	84	97	71	79
Infinitive clause	85	66	55	89	92	87	80	88	67	95	95
Object clause	16	56	44	94	69	62	75	100	88	81	81
Pseudo-cleft sentence	73	90	86	68	73	85	89	92	84	77	74

Accuracy of the metrics(%) with regards to the linguistically-motivated phenomena for English-German. Continued on next page

727

QE as a metric	baselines		ref. based metrics																		
			CoMeT-QE-Ensemble									CoMeT-QE									
#	BERTScore	BLEURT-20	BLEU	chrf	yisf-1	COMET	CoMeT-QE	CoMeT-QE-XXL	GEMBA-MoM	KG-BERT	MS-CoMeT-QE	MS-CoMeT-QE	Metrix-23-QE-b	Metrix-23-QE-c	Metrix-23-QE-E	Metrix-23-QE-F	Metrix-23-QE-G	Metrix-23-QE-H	MS-CoMeT-QE-Ensemble	CoMeT-QE-Ensemble	
ling. phenomenon	Reflex. - future II progr.	Reflex. - future II simple	Reflex. - past perf. progr.	Reflex. - past perf. simple	Reflex. - past pres. progr.	Reflex. - past pres. simple	Reflex. - simple past	Reflex. - simple pres.	Reflex. - simple trans.	Reflex. - simple trans. - cond. I progr.	Reflex. - simple trans. - cond. I simple	Reflex. - simple trans. - cond. II progr.	Reflex. - simple trans. - cond. II simple	Reflex. - simple trans. - cond. II pres.	Reflex. - simple trans. - cond. II pres. progr.	Reflex. - simple trans. - cond. II simple	Reflex. - simple trans. - cond. II simple progr.	Reflex. - simple trans. - cond. II simple pres.	Reflex. - simple trans. - cond. II simple pres. progr.		
	81	65	64	64	75	80	68	57	68	75	58	85	77	85	52	85	89	86	51	83	37
	56	71	66	77	54	88	88	71	64	77	95	89	84	73	84	39	98	98	50	100	49
	98	60	50	67	64	71	66	49	51	81	79	70	88	85	88	52	74	73	83	39	77
	53	62	47	50	60	100	100	100	0	60	40	40	100	100	100	40	40	40	60	100	100
	33	76	48	88	67	76	76	45	48	61	100	88	85	88	48	100	100	100	40	60	100
	39	59	46	67	59	69	72	46	44	97	72	69	64	69	31	74	85	85	38	82	62
	99	62	51	54	56	56	56	37	62	51	75	66	46	91	46	39	58	76	70	41	79
	119	71	65	73	73	73	73	77	64	71	63	73	82	82	45	95	98	100	34	71	55
	138	65	65	67	67	88	63	64	67	55	70	75	58	87	58	47	97	73	56	72	78
	11	73	82	73	91	64	82	82	82	100	73	100	100	91	100	100	100	100	91	100	100
	11	55	91	45	64	36	82	52	55	91	36	91	100	55	36	55	82	91	45	100	100
	9	67	100	89	89	56	100	22	100	67	89	100	56	56	89	100	47	89	33	89	44
	20	70	55	75	65	80	55	85	60	70	95	95	70	75	70	40	75	85	40	65	65
	2	50	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100
	12	42	83	75	83	25	50	75	75	75	75	75	58	92	67	92	92	58	75	75	92
	22	68	95	64	68	68	77	50	95	73	86	36	77	57	86	57	73	64	68	68	73
	33	62	59	77	67	87	87	62	80	80	80	80	80	80	80	80	80	80	80	80	80
	16	50	69	50	81	81	56	69	94	94	38	94	31	94	75	81	50	75	75	81	38
	9	44	78	89	89	33	89	22	78	100	89	100	100	44	100	56	100	78	89	67	100
	5	20	80	80	80	20	80	44	78	44	80	100	100	22	100	100	22	100	100	20	20
	9	33	67	78	78	44	78	44	78	100	89	100	100	44	100	33	67	78	100	22	22
	10	30	70	20	40	30	40	30	50	30	40	50	30	40	50	50	30	40	40	20	40
	23	61	43	96	70	35	57	87	52	91	78	91	96	91	96	83	48	96	78	30	67
	16	31	62	38	75	69	62	56	94	50	81	100	62	100	100	94	75	50	81	94	100
Verb valency																					
Case government																					
Catenative verb																					
Middle voice																					
Passive voice																					
Resultative																					
macro avg.	8945	67	65	74	75	70	71	73	67	73	76	77	68	68	84	80	83	43	82	52	71
micro avg.	8945	70	65	75	76	72	69	74	68	73	74	75	74	68	74	60	62	74	71	77	71

Table 5: Accuracy of the metrics(%) with regards to the linguistically-motivated phenomena for English-German

Table 6: Accuracy of the metrics (%) with regards to the linguistically-motivated phenomena for English-Russian