

# Machine Translation Metrics are better in evaluating Linguistic Errors on LLMs than on Encoder-Decoder Systems

Eleftherios Avramidis, Shushen Manakhimova, Vivien Macketanz and Sebastian Möller

German Research Center for Artificial Intelligence (DFKI),

Berlin, Germany

firstname.lastname@dfki.de

## Abstract

This year’s MT metrics challenge set submission by DFKI expands previous years’ linguistically motivated challenge set. It includes 137,000 items extracted from 100 MT systems for the two language directions (en→de, en→ru), covering more than 100 linguistically motivated phenomena organized in 14 linguistic categories. The metrics with the statistically significant best performance with regard to our linguistically motivated analysis are METRICX-24-HYBRID and METRICX-24 for en→de and METRICX-24 for en→ru, whereas METAMETRICS and XCOMET are in the next ranking positions in both language pairs. Metrics are more accurate in detecting linguistic errors among LLM translations than in translations based on the encoder-decoder NMT architecture. Some of the most difficult phenomena for the metrics to score are the transitive past progressive, the multiple connectors, and the ditransitive simple future I for en→de and the pseudogapping, the contact clause and the cleft sentences for en→ru. Despite its overall low performance, the LLM-based metric GEMBA performs best in scoring German negation errors.

## 1 Introduction

For almost two decades, the development and evaluation of machine translation (MT) have relied on automatic metrics. MT metrics aim to digest and automate various aspects of human judgment of MT output into numerical scores. Over the years, these metrics have undergone several technological changes (from measuring overlap to grammatical features and neural models). Still, at the same time, they have had to follow the technological evolution of MT systems, moving from phrase-based statistical systems to NMT encoder-decoder models and, more recently, to large language models (LLMs). As we witness the first efforts to use and evaluate LLMs in the task of MT, it is of great interest to see to what extent pre-existing MT methodologies

can adapt to the needs of the new technologies. An obvious question is to what extent MT metrics developed and tested for NMT can be applied to evaluating LLMs.

This year’s Metrics Task (WMT24; Freitag et al., 2024) provides a very good opportunity to evaluate the metrics under these particular circumstances, as the evaluated MT outputs have for the first time been produced by numerous LLMs (Kocmi et al., 2024). Meanwhile, the ability of LLMs to act as judges for translations is being explored through the participation of an LLM-based metric.

Given this perspective, this paper extends previous work on linguistically motivated challenge sets for MT metrics to investigate whether LLMs can influence MT evaluation. As part of this year’s submission to the challenge set subtask of the WMT24 Metrics Task, we repeat the methodology of previous years to evaluate the metrics on a controlled test set that can rank them with regard to their ability to detect linguistic errors by providing fine-grained statistics for each linguistic phenomenon. We then analyze whether the metrics perform differently on MT output from LLMs as opposed to output from encoder-decoder systems. In addition, we see in which linguistic aspects the LLM-based metric performs better or worse than the specialized metrics.

The rest of the paper is structured as following: Section 2 describes briefly the generation of the challenge set. Section 3 presents and discusses the results, whereas the conclusion is given in section 4

## 2 Method

This year’s linguistically-motivated challenge set is an extension of the challenge sets that were submitted the previous years (Avramidis and Macketanz, 2022; Avramidis et al., 2023).

The source sentences  $s$  originate from an MT evaluation test suite (Macketanz et al., 2022a). Each sentence has been carefully constructed to test one particular phenomenon. Every phenomenon is

tested by more sentences (with a minimum of 20 sentences), whereas the phenomena are aggregated in a few categories. At the moment, there are more than 100 phenomena and 14 categories.

As part of the WMT shared tasks of the previous years, these source sentences have been given to a large amount of MT systems, and their output has been evaluated by combining regular expressions and annotations by linguists, labeling every output as correct ( $t \in T$ ) or incorrect ( $\hat{t} \in T'$ ).

In order to use this test set to evaluate the MT metrics, we create examples in the form of  $(s, \hat{t}, t, r) \in S$ , where each example contains one source sentence  $s$ , one incorrect translation hypothesis  $\hat{t}$ , one correct translation hypothesis  $t$  and one reference translation  $r$ . The correct translation hypotheses  $t$  and the reference translations  $r$  are sampled with permutations from the same set of correct translations  $T$ . Then, we decompose the set of examples  $S$  into a blind test set  $S'$ , where each example includes either an incorrect translation  $(s, \hat{t}, r)$  or a correct translation  $(s, t, r)$  along with the source and the reference. The separated contrastive examples are shuffled, and we set aside a file that contains the golden truth, indicating which samples are correct or incorrect.

As part of the Metrics Task, every shuffled translation  $t$  and  $\hat{t}$  is scored by every  $M$ , given the reference  $r$  in the given blind test set  $S'$ , without knowing if it is correct or incorrect. A contrastive pair scoring is considered correct if the metric delivers a score for the incorrect translation hypothesis, which is lower than the one of the correct translation hypothesis  $M(s, \hat{t}, r) < M(s, t, r)$ . Finally, for every phenomenon and category and for every metric, the respective accuracy is calculated by dividing the number of correctly scored contrastive pairs by the total amount of examples.

$$\text{acc}_M = \frac{|M(s, \hat{t}, r) < M(s, t, r)|}{|(s, \hat{t}, t, r)|}$$

$$(s, \hat{t}, r) \cup (s, t, r) \in S' \quad (s, \hat{t}, t, r) \in S$$

Lastly, we provide three types of score averaging:

- i) **Micro-average:** This approach treats all items equally, aggregating all test items to compute the average percentages.
- ii) **Category macro-average:** Here, all categories are treated equally, with the percent-

ages 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.

The current version of the challenge set contains MT outputs from the WMT Shared Tasks of the years 2019-2024 (Avramidis et al., 2019, 2020; Macketanz et al., 2021, 2022b; Manakhimova et al., 2023, 2024). The English to German version contains 39,463 contrastive pairs, while the English to Russian version contains 30,108 pairs.

## 3 Results

### 3.1 English-German

The comparison of the metrics based on the accuracies per category for English-German can be seen in table 2, whereas the detailed phenomena in table 4. One can see that the metrics which have the highest accuracy with statistical significance are METRICX24-HYBRID and METRICX24 (Juraska et al., 2024), with more than 80.7 % macro-average. Both metrics are very good at multi-word expressions (mostly verbal MWEs). The former is the best of all metrics at coordination/ellipsis and non-verbal agreement (genitive and personal pronoun coreference). In contrast, the latter performs best at verb valency (resultative and passive voice). The metrics ‘‘METAMETRICS’’ (Anugraha et al., 2024) and XCOMET (Guerreiro et al., 2023) follow in the ranking, with more than 80% macro-averaged accuracy.

The LLM-based metric GEMBA (Kocmi and Federmann, 2023) performs relatively low, with an average accuracy of 69.7%, even below the baseline non-tuned metric CHRf (Popović, 2015). It is nevertheless remarkable that this metric has the best score on negation, among all metrics (97.4%, 4.5% higher than the best system). The fact that most of the metrics will miss 10% of the negations is rather noteworthy, given the implications of such a mistake on the meaning of the sentence. It is also remarkable that a reference-less metric, METRICX24-HYBRID-QE, achieves the highest accuracy on long-distance dependencies and interrogatives, mainly on the phenomenon of negative inversion.

Some of the most difficult phenomena for the

	METRICX24	METRICX24-HYB	METAMETRICS	XCOMET
encdec vs. encdec	73.2	72.3	70.8	69.7
LLM vs. encdec	77.3	76.9	79.9	77.6
LLM vs. LLM	79.9	78.1	80.0	79.1

Table 1: Accuracy of the metrics when they evaluate contrastive pairs containing (a) MT output only by encoder/decoder systems, (b) one encoder/decoder output and one LLM output, (c) only LLM output

metrics to score are transitive past progressive, multiple connectors, and ditransitive simple future I.

### 3.2 English-Russian

The comparison of the metrics based on the accuracies per category for English-Russian can be seen in table 3, whereas the detailed phenomena in table 5. MetricX-24 is the clear winner in this language direction, achieving a macro-averaged accuracy of 82.5% MetricX-24 excels in ambiguity, false friends, non-verbal agreement (coreference & genitive), verb semantics, and verb valency. The ranking of the metrics is similar to the one for English-German, with METAMETRICS, METRICX24-HYBRID and XCOMET having the next position, with more than 79.6% accuracy in macro-average.

If one focuses again on the phenomenon of negation, they would notice that in English-Russian, the highest accuracy is achieved by the baseline metric CHRf, whereas most metrics perform here very low (61% on average) Some of the most difficult phenomena for this language direction are the pseudogapping, the contract clause, and the cleft sentences for en→ru.

### 3.3 Comparing performance of metrics over LLM vs. encoder-decoder systems

Table 1 presents the accuracies of the 4 best performing metrics on three subsets of the challenge sets. Here every subset contains contrastive pairs which consist of

- (a) two MT outputs, both by encoder/decoder NMT systems
- (b) one encoder/decoder and one LLM output
- (c) two LLM outputs

One can see that all four metrics exhibit higher accuracy when scoring contrastive translations originating from LLMs. This indicates that despite the fact that LLM translations achieve very good performance (Kocmi et al., 2024), their fewer errors are easier to be distinguished by the automatic metrics. Whether there is a systematic reason for this phenomenon remains to be investigated.

## 4 Conclusion

We presented the MT metrics challenge set of DFKI for two language directions (en-de, en-ru). This year, we have expanded the set to include outputs from encoder-decoder NMT systems and LLMs. The number of test items (total of 137,000) allows for producing fine-grained scores for every linguistic phenomenon and statistically significant comparisons among the MT metrics. We also identified the best-performing metric, METRICX-24, for both language directions.

### Acknowledgements

This research was supported by the German Research Foundation (Deutsche Forschungsgemeinschaft; DFG) through the project TextQ. 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 to the creation of the test suite.

### References

- David Anugraha, Garry Kuwanto, Lucky Susanto, Derry Tanti Wijaya, and Genta Indra Winata. 2024. Metametrics-MT: Tuning machine translation metametrics via human preference calibration. In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.
- 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, 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.
- Eleftherios Avramidis, Shushen Manakhimova, Vivien Macketanz, and Sebastian Möller. 2023. [Challenging the state-of-the-art machine translation metrics from a linguistic perspective](#). In *Proceedings of the Eighth Conference on Machine Translation*, pages 713–729, Singapore. Association for Computational Linguistics.
- Markus Freitag, Nitika Mathur, Daniel Deutsch, Chikui Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva, and Alon Lavie. 2024. [Are llms breaking mt metrics? results of the wmt24 metrics shared task](#). In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.
- Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2023. [xcomet: Transparent machine translation evaluation through fine-grained error detection](#).
- Juraj Juraska, Daniel Deutsch, Mara Finkelstein, and Markus Freitag. 2024. [Metricx-24: The google submission to the WMT 2024 metrics shared task](#). In *Proceedings of the Ninth Conference on Machine Translation*, 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, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinþór Steingrímsson, and Vilém Zouhar. 2024. [Findings of the WMT24 general machine translation shared task: the LLM era is here but mt is not solved yet](#). In *Proceedings of the Ninth Conference on Machine Translation*, USA. Association for Computational Linguistics.
- Tom Kocmi and Christian Federmann. 2023. [GEMBA-MQM: Detecting translation quality error spans with GPT-4](#). In *Proceedings of the Eighth Conference on Machine Translation*, pages 768–775, Singapore. 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.
- Shushen Manakhimova, Eleftherios Avramidis, Vivien Macketanz, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian Möller. 2023. [Linguistically motivated evaluation of the 2023 state-of-the-art machine translation: Can ChatGPT outperform NMT?](#) In *Proceedings of the Eighth Conference on Machine Translation*, pages 224–245, Singapore. Association for Computational Linguistics.
- Shushen Manakhimova, Vivien Macketanz, Eleftherios Avramidis, Ekaterina Lapshinova-Koltunski, Sergei Bagdasarov, and Sebastian Möller. 2024. [Investigating the linguistic performance of large language models in machine translation](#). In *Proceedings of the Ninth Conference on Machine Translation*, USA. 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*, pages 392–395.

## A Accuracies per category

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

ling. category	#	metric																				avg						
		MetricX-24-Hybrid	MetricX-24	metametrics	XCOMET	BLEURT-20	COMET-22	CometKwi-XXL	MetricX-24-QE	MetricX-24-Hybrid-QE	XCOMET-QE	YISI-1	sentinel-cand-mqm	CometKwi	MEE4	chrtS	BERTScore	chrf	gamba	spBLEU	damononli		mononli	BLEU	XLsimMqm	XLsimDA	PrismRefSmall	PrismRefMedium
Ambiguity	4614.0	85.1	85.9	<b>89.9</b>	80.0	<b>89.7</b>	<b>89.5</b>	60.8	74.6	70.6	61.9	88.6	77.7	48.2	82.1	83.8	78.2	85.2	70.0	80.0	83.8	83.0	64.1	55.2	55.2	68.1	60.6	75.1
Coordination & ellipsis	4373.0	<b>81.3</b>	74.2	74.4	77.4	76.5	76.7	80.2	78.2	78.8	74.4	69.2	76.7	71.1	63.8	62.9	67.3	62.2	66.5	61.8	61.0	62.9	60.6	49.5	49.5	51.1	49.1	67.6
False friends	1389.0	79.9	78.2	78.3	73.9	72.7	<b>85.9</b>	85.2	69.8	74.3	71.1	73.1	77.0	72.1	80.4	77.1	69.2	74.9	81.9	65.9	48.6	38.2	64.1	78.3	74.3	74.3	58.8	72.4
Function word	1900.0	78.1	80.6	82.2	86.0	81.9	<b>87.3</b>	83.0	81.7	78.6	82.2	72.9	85.8	86.7	76.9	77.2	<b>86.9</b>	74.4	70.9	74.2	64.9	60.1	78.6	55.7	55.7	52.2	53.8	74.9
LDD & interrogatives	1002.0	83.4	80.1	80.8	80.6	80.3	74.5	78.7	81.8	<b>84.7</b>	81.7	59.1	68.6	78.4	64.2	59.3	64.1	57.5	58.6	61.5	66.7	64.5	60.5	62.0	62.0	49.6	47.7	68.9
MWE	5816.0	<b>87.0</b>	<b>87.3</b>	85.9	86.2	84.1	82.9	80.0	82.5	81.2	80.5	80.3	84.0	76.4	75.1	75.9	76.1	73.3	82.0	70.7	77.3	76.6	71.5	67.0	67.0	59.4	55.6	77.1
Named entity & terminology	22891.0	71.5	74.2	74.2	68.8	71.7	73.6	58.0	55.3	60.9	56.9	74.7	52.1	50.2	72.2	70.5	67.1	68.9	48.1	70.0	<b>75.4</b>	73.1	62.0	48.5	48.5	49.8	50.1	63.3
Negation	506.0	92.9	89.5	88.5	91.1	92.7	92.9	93.3	93.9	91.3	90.9	87.9	74.5	95.3	90.7	82.8	86.0	76.7	<b>97.4</b>	73.9	86.6	88.3	73.7	58.3	58.3	58.5	58.1	83.2
Non-verbal agreement	15497.0	<b>83.6</b>	80.6	77.4	80.9	78.2	73.3	80.2	82.3	82.4	79.2	65.7	76.2	72.9	65.6	66.1	63.7	65.9	72.7	64.3	59.6	59.5	62.2	57.8	57.8	51.0	49.0	69.5
Punctuation	2435.0	62.2	64.4	64.9	63.2	71.9	72.4	70.1	70.4	65.9	64.9	71.6	<b>80.1</b>	71.3	69.9	72.1	66.0	68.5	44.3	67.3	66.9	50.7	68.7	50.3	50.3	50.6	50.8	64.2
Subordination	4698.0	89.1	87.5	86.3	<b>89.3</b>	84.1	83.9	89.2	<b>89.5</b>	<b>89.4</b>	86.9	78.9	80.8	<b>89.8</b>	76.6	76.1	76.4	74.1	72.6	72.3	66.1	70.9	73.9	44.4	44.4	57.5	54.3	76.3
Verb tense/aspect/mood	10120.0	78.6	81.8	79.2	<b>83.4</b>	73.0	68.2	80.6	77.4	77.3	81.8	67.5	52.2	72.1	67.8	67.8	65.9	66.7	73.3	63.0	63.6	66.0	62.6	51.8	51.8	59.1	52.7	68.7
Verb valency	3486.0	80.8	<b>84.6</b>	<b>84.5</b>	81.7	81.7	77.0	83.8	82.9	80.3	80.9	73.7	75.5	73.2	67.4	71.1	67.9	71.6	67.4	67.2	70.5	71.3	62.3	61.2	61.2	55.0	53.9	72.6
macro avg.	78727.0	<b>81.0</b>	<b>80.7</b>	80.5	80.2	79.9	79.9	78.7	78.5	78.1	76.4	74.1	73.9	73.7	73.3	72.5	71.9	70.8	69.7	68.6	68.5	66.5	66.5	56.9	56.9	56.6	53.4	71.8
micro avg.	78727.0	<b>79.1</b>	<b>79.4</b>	78.6	77.9	77.1	76.0	73.3	73.1	74.2	72.0	72.8	67.3	66.4	70.9	70.7	68.8	69.5	64.8	68.0	68.8	67.9	64.3	53.9	53.9	54.4	51.9	69.0

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

ling. category	#	metric																				avg					
		MetricX-24	metametrics	MetricX-24-Hybrid	XCOMET	BLEURT-20	COMET-22	CometKiwI-XXL	MetricX-24-QE	XCOMET-QE	MetricX-24-Hybrid-QE	YISI-1	CometKiwI	BERTscore	sentinel-cand-mqnm	chrf	chrf	spBLEU	BLEU	dammomni	monnomli		gemba	XLsimMqnm	XLsimDA	PrismRefSmall	PrismRefMedium
Ambiguity	3788.0	96.4	95.1	93.2	89.8	87.4	80.9	96.3	83.8	91.4	82.6	77.2	75.3	87.1	74.7	73.1	70.6	68.9	80.5	78.4	89.4	89.4	43.9	43.9	48.1	45.3	78.0
Coordination & ellipsis	2273.0	80.6	81.5	80.4	74.9	76.6	81.0	81.4	77.2	81.8	68.6	78.6	68.1	75.5	63.5	61.5	62.7	62.7	65.4	66.3	60.1	52.5	52.5	47.7	48.7	69.2	
False friends	2414.0	87.8	86.3	76.3	82.4	83.0	69.1	69.2	68.6	68.4	87.7	52.4	76.3	58.0	84.9	83.2	80.7	75.8	80.8	62.2	34.0	43.2	43.2	53.1	42.0	69.3	
Function word	2433.0	82.5	78.0	73.4	79.7	81.4	83.0	85.7	86.3	71.8	65.7	79.3	69.3	82.7	64.8	60.3	65.6	73.2	56.4	57.0	56.3	73.7	73.7	50.3	49.0	71.3	
LDD & interrogatives	1939.0	85.4	86.0	87.8	84.8	81.8	82.6	84.9	87.3	83.4	87.6	65.5	77.6	66.2	87.8	62.5	59.9	62.0	61.5	55.1	58.3	68.1	54.2	54.2	51.6	46.4	71.3
MWE	9602.0	82.9	82.9	81.2	80.5	81.4	82.2	81.6	77.7	83.9	77.0	75.6	74.6	72.3	75.5	73.1	72.8	70.9	67.9	65.1	69.5	53.5	53.5	51.7	51.0	72.8	
Named entity & terminology	16284.0	82.8	84.9	81.6	80.6	84.9	84.3	71.6	72.2	69.5	71.6	83.0	71.3	78.8	70.6	80.3	78.7	78.3	72.5	72.9	67.7	64.1	47.6	47.6	53.6	52.1	72.1
Negation	346.0	65.3	59.8	58.7	49.4	72.3	67.3	58.7	57.8	45.4	44.5	79.5	49.4	80.3	41.3	82.9	83.5	74.3	72.3	70.5	72.5	42.5	49.7	49.7	49.4	45.7	60.9
Non-verbal agreement	6755.0	86.4	81.5	84.4	82.3	79.6	77.4	78.7	83.0	80.9	81.5	72.1	77.6	68.7	73.1	69.4	68.2	67.4	64.6	59.6	60.9	68.4	68.4	68.4	51.2	47.5	72.1
Punctuation	363.0	73.0	71.1	72.7	71.3	68.6	76.0	75.8	63.6	70.8	67.2	75.8	72.2	73.3	70.5	62.0	58.4	64.7	60.9	51.0	64.5	60.9	57.3	57.3	46.3	45.5	65.2
Subordination	6625.0	74.7	74.5	71.4	75.0	72.7	77.1	78.4	72.5	75.2	71.7	69.3	68.6	66.9	73.5	63.8	62.4	64.3	64.0	56.4	63.3	53.6	50.5	50.5	51.0	48.1	66.0
Verb semantics	275.0	88.0	82.2	88.0	85.5	86.5	74.2	79.6	75.3	80.0	76.0	53.1	69.8	55.6	55.3	60.7	65.1	53.8	48.7	68.4	66.5	72.0	33.5	33.5	65.5	67.3	67.4
Verb tense/aspect/mood	2994.0	85.0	86.0	82.6	86.2	75.5	79.7	85.8	82.8	80.7	79.3	69.7	72.6	68.7	70.4	68.1	66.7	68.8	63.1	60.1	55.9	61.6	47.5	47.5	50.6	51.3	69.9
Verb valency	3022.0	83.3	82.3	80.4	83.5	76.8	76.3	80.9	82.0	81.6	81.5	69.6	73.8	72.8	72.0	72.2	71.8	69.0	67.7	66.6	64.2	66.9	60.7	60.7	51.6	46.7	71.8
macro avg.	59113.0	82.5	80.6	80.5	79.6	79.0	78.9	77.9	77.9	75.8	75.6	72.8	71.1	71.1	70.7	70.4	69.0	68.2	66.2	65.1	64.5	62.0	52.6	52.6	51.6	49.0	69.8
micro avg.	59113.0	83.4	82.8	81.7	81.4	80.7	81.0	77.8	78.7	76.2	77.5	75.8	72.9	72.9	73.0	73.1	71.4	71.2	68.5	66.8	64.8	64.4	52.8	52.8	51.7	49.3	71.3









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

ling. category	ling. phenomenon	#	metric																									
			XCOMET	MetricX-24	metametrics	MetricX-24-QE	MetricX-24-Hybrid-QE	XCOMET-QE	CometKwi-XXL	BLEURT-20	CometKwi	COMET-22	chrtS	MEE4	chrt	gemba	BERTScore	YISI-1	spBLEU	BLEU	monnonih	dannonnih	sentinel-cand-mqm	PrismRetSmall	XLsimMqm	XLsimDA	XLsimMqm	PrismRetMedium
Verb valency	Transitive - simple present	35	80	71	63	74	74	80	43	94	51	37	37	46	63	29	34	43	49	49	57	83	71	66	66	99	63	61
	Case government	189	81	76	78	87	78	77	79	78	77	69	68	62	81	70	77	67	62	60	73	70	51	51	51	51	52	72
	Catenaive verb	885	89	88	81	87	79	86	81	70	74	65	64	67	76	63	68	60	50	50	56	60	56	73	73	73	53	72
	Mediopassive voice	183	95	95	99	99	94	96	91	92	97	96	96	98	87	93	95	93	89	79	80	89	82	63	63	75	89	
	Passive voice	176	77	84	81	84	81	82	65	62	67	83	78	55	47	52	74	51	44	44	44	57	78	51	47	47	55	64
Verb semantics	Resultative	1203	83	88	85	83	85	82	86	80	76	68	76	76	67	67	74	69	67	75	80	85	49	58	58	58	76	
	Semantic roles	670	65	77	55	87	72	71	79	88	49	74	74	71	34	76	77	73	70	80	79	81	58	55	55	55	41	69
	Verb semantics	180	87	73	69	82	73	71	78	68	62	61	53	57	55	69	58	55	53	56	50	54	59	69	69	69	54	64
	macro avg.	78727	82	82	81	79	78	77	77	77	73	71	71	70	69	69	68	67	65	65	64	63	59	57	57	57	57	70
	micro avg.	78727	78	79	79	73	74	72	73	77	76	76	71	71	70	65	69	73	68	64	68	69	67	54	54	54	54	52

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

ling. category	ling. phenomenon	#	metric																									
			MetricX-24	metametrics	XCOMET	CometKwi-XXL	MetricX-24-Hybrid	MetricX-24-QE	COMET-22	BLEURT-20	MetricX-24-Hybrid-QE	XCOMET-QE	sentinel-cand-mqm	CometKwi	YISI-1	BERTScore	chrtS	spBLEU	chrt	BLEU	monnonih	dannonnih	gemba	XLsimDA	XLsimMqm	PrismRetSmall	PrismRetMedium	avg
Ambiguity	Lexical ambiguity	3788	97	96	93	81	95	96	63	67	75	75	75	77	83	75	71	73	69	78	81	89	44	44	48	45	78	
	Gapping	698	93	92	93	87	92	91	93	88	89	90	85	96	86	85	81	81	77	79	81	74	78	57	57	55	56	81
	Pseudogapping	381	71	67	55	64	70	67	62	63	70	54	54	57	60	59	57	54	57	55	55	53	39	39	40	40	57	
	Right node raising	183	78	74	74	77	77	70	75	70	70	70	70	79	66	68	64	60	61	58	70	63	55	45	45	45	61	67
	Sluicing	384	80	82	86	89	84	82	82	77	73	82	85	84	77	61	63	56	54	62	57	54	48	48	48	35	41	67
False friends	Stripping	375	70	69	80	76	74	79	63	67	81	75	69	69	62	60	56	53	54	60	68	78	35	57	57	49	50	
	VP-ellipsis	252	80	77	80	88	85	86	75	75	90	73	76	83	56	56	48	53	49	49	45	55	70	65	65	44	66	
	False friends	2414	88	84	76	69	86	69	83	82	68	69	58	88	76	85	81	83	76	62	81	34	43	43	53	42	69	
	Focus particle	846	70	62	63	68	60	67	67	57	49	66	63	55	65	63	64	66	60	67	57	50	44	71	71	45	50	
	Question tag	1587	89	87	95	91	81	95	89	92	84	97	78	85	92	66	73	65	61	77	57	60	63	75	75	53	48	77
LDD & interrogatives	Inversion	333	82	88	84	73	90	89	79	83	92	78	85	67	71	68	68	66	68	67	61	60	63	47	47	56	52	71
	Modifying Comparison	90	68	71	74	87	69	78	52	100	71	74	98	56	44	41	33	29	29	28	56	67	67	73	73	37	40	
	Multiple connectors	400	97	92	93	95	98	96	88	78	96	96	95	94	64	67	62	58	61	55	60	52	86	52	52	69	52	

Continued on next page



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

ling. category	ling. phenomenon	#	metric																				avg						
			MetricX-24	metametrics	XCOMET	CometKwi-XXL	MetricX-24-Hybrid	MetricX-24-QE	COMET-22	BLEURT-20	MetricX-24-Hybrid-QE	XCOMET-QE	sentinel-cand-mqm	CometKwi	YSI-1	BERTScore	chrs	sBLEU	chrf	BLEU	monnoli	dannnoli		gemba	XLsimDA	XLsimMqm	PrismRetSmall	PrismRetMedium	
Verb valency	Imperative	575	88	<b>89</b>	85	<b>88</b>	84	87	83	79	84	85	80	68	74	74	66	69	64	63	69	69	50	44	44	42	44	71	
	Intransitive	103	94	91	87	<b>98</b>	94	90	94	81	94	88	93	93	81	68	68	65	65	54	58	60	56	42	42	37	41	73	
	Reflexive	514	88	89	84	<b>94</b>	77	90	84	77	77	83	85	67	70	69	70	70	68	65	41	55	88	58	58	51	51	72	
	Transitive	516	78	<b>87</b>	85	85	80	80	77	68	80	66	47	51	77	74	73	75	78	69	50	61	28	50	50	63	59	68	
	Case government	331	76	<b>86</b>	78	71	76	74	77	85	75	71	51	70	75	82	83	75	82	71	81	79	78	74	74	43	40	73	
	Catenative verb	358	72	73	69	68	<b>73</b>	71	66	70	72	70	67	63	58	63	66	61	62	59	62	63	65	49	49	50	48	64	
	Impersonal Subject	217	86	77	82	95	85	<b>98</b>	75	71	90	94	74	87	65	67	61	53	60	53	76	75	59	50	50	43	41	71	
	Mediopassive voice	409	77	79	89	85	69	<b>89</b>	80	80	72	77	<b>90</b>	83	82	73	75	65	70	63	67	57	63	64	58	38	37	70	
	Passive voice	228	94	89	87	84	92	<b>98</b>	88	84	80	84	69	82	79	83	75	75	74	84	73	83	66	74	74	52	44	79	
	Resultative	660	<b>91</b>	86	<b>91</b>	88	78	80	88	87	88	87	82	66	72	75	73	76	73	65	67	73	62	62	62	64	55	76	
	Semantic roles	270	<b>91</b>	82	83	77	<b>85</b>	73	79	89	87	79	60	80	71	67	79	70	81	65	80	84	53	63	63	63	57	50	74
	Verb semantics/Verb semantics	549	81	<b>83</b>	83	80	74	73	71	70	78	78	74	75	66	70	72	68	72	66	43	46	66	58	58	53	48	68	
	macro avg.		59113	<b>82</b>	81	81	81	81	80	79	79	78	75	71	70	68	67	66	64	64	63	63	62	54	54	50	48	70	
micro avg.		59113	<b>83</b>	83	81	78	82	79	81	81	77	76	73	73	76	73	71	71	68	65	67	64	53	53	52	49	71		