Guardians of the Machine Translation Meta-Evaluation Sentinel Metrics Fall In!

Anonymous ACL submission

Abstract

Annually, the organizers of the Metrics Shared Task at the Conference on Machine Translation (WMT) conduct the meta-evaluation of Machine Translation (MT) metrics, ranking them according to their correlation with human judgments. Their results guide researchers toward enhancing the next generation of metrics and MT systems. With the recent introduction of neural metrics, the field has witnessed notable advancements. Nevertheless, the inherent opacity of these metrics has posed substantial challenges to the meta-evaluation process. This work highlights two issues with the metaevaluation framework currently employed in WMT, and assesses their impact on the metrics rankings. To do this, we introduce the concept of sentinel metrics, which are designed explicitly to scrutinize the accuracy, robustness, and fairness of the meta-evaluation process. By employing sentinel metrics, we aim to validate our findings, and shed light and monitor the potential biases or inconsistencies in the rankings. We discover that the present meta-evaluation framework favors two categories of metrics: i) those explicitly trained to mimic human quality assessments, and ii) continuous metrics. Ultimately, we raise concerns regarding the evaluation capabilities of state-of-the-art metrics, highlighting that they might be basing their assessments on spurious correlations found in their training data.

1 Introduction

Over the past few years, the Machine Translation (MT) field has witnessed significant advancements, largely driven by the advent of neural architectures, with the Transformer (Vaswani et al., 2017) being the most notable. Modern MT systems now deliver translations that are mostly fluent and accurate, posing a challenge for their quality evaluation – even when conducted by human annotators, especially those who lack professional training (Freitag et al., 2021a). Under these circumstances, shallow

overlap-based metrics are gradually being replaced by neural-based metrics, that demonstrate a better correlation with human judgments (Freitag et al., 2022). However, a significant limitation is that most neural metrics are black-box systems trained to predict human judgments in the form of scalar scores, and typically do not provide justifications for their assessments. Besides rendering them challenging to interpret, such opacity also complicates their meta-evaluation. In this respect, we found that certain strategies for the assessment of MT metrics' capabilities - which have recently been employed in the context of the Metrics Shared Task at the Conference on Machine Translation (WMT)¹ - favor specific metric categories and potentially encourage undesirable metrics behavior. To demonstrate these problems, we introduce the concept of sentinel metrics, i.e., a suite of metrics serving as a probe to identify pitfalls in the meta-evaluation process. Sentinel metrics are either trained with incomplete information - which makes them inherently unable to properly evaluate the quality of machine-translated text - or consist of variations of existing metrics - which have been devised to expose specific issues in the meta-evaluation.

As an example, in Table 1, we present the segment-level ranking of WMT23 with the inclusion of a sentinel metric. As can be seen, SENTINEL_{CAND} ranks in the upper half. SENTINEL_{CAND} is a sentinel metric designed to assess the quality of a candidate translation solely based on the translation itself, without accessing its source sentence or any reference translation. Arguably, such a metric should only be capable of evaluating a translation's fluency, but not its adequacy in conveying the original message, and a fair assessment should rank it at lower positions. Notably, SENTINEL_{CAND} is above strong baselines

¹With its first edition in 2006 (Koehn and Monz, 2006), "WMT is the main event for machine translation and machine translation research." (https://machinetranslate.org/wmt).

Metric		Avg. corr
XCOMET-Ensemble	1	0.697
MetricX-23	2	0.682
XCOMET-QE-Ensemble*	3	0.681
MetricX-23-QE*	4	0.681
mbr-metricx-qe*	5	0.652
GEMBA-MQM*	6	0.639
MaTESe	7	0.636
CometKiwi*	8	0.632
sescoreX	9	0.628
SENTINEL _{CAND} *	10	0.626
cometoid22-wmt22*	11	0.625
KG-BERTScore*	12	0.624
COMET	13	0.622
BLEURT-20	14	0.622
Calibri-COMET22-QE*	15	0.603
Calibri-COMET22	16	0.603
<u>YiSi-1</u>	17	0.600
docWMT22CometDA	18	0.598
docWMT22CometKiwiDA*	19	0.598
prismRef	20	0.593
MS-COMET-QE-22*	21	0.588
BERTscore	22	0.582
mre-score-labse-regular	23	0.558
XLsim	24	0.544
f200spBLEU	25	0.540
MEE4	26	0.539
tokengram_F	27	0.537
chrF	28	0.537
BLEU	29	0.533
prismSrc*	30	0.530
embed_llama	31	0.529
eBLEU	32	0.491
Random-sysname*	33	0.463

Table 1: Segment-level ranking of the primary submissions to the WMT 2023 Metrics Shared Task, with the inclusion of sentinel metrics. The values in the column 'Avg. corr' are obtained by averaging the correlations of the 6 segment-level tasks of WMT 2023. Starred metrics are reference-free, underlined metrics are baselines, and highlighted metrics are sentinels. Ranks represent clusters of statistical significance and are computed following Freitag et al. (2023), which leverage the PERM-BOTH hypothesis test introduced by Deutsch et al. (2021). In Table 3 in the Appendix, we report the metrics performance in terms of rank and correlation in all the 6 tasks that contribute to this ranking. All the rankings present in this work have been computed with the official shared task library (https://github.com/ google-research/mt-metrics-eval).

such as COMET (Rei et al., 2020) and BLEURT-20 (Sellam et al., 2020), suggesting that there might be some issues with the segment-level meta-evaluation methods used in WMT23.

In this work, we: i) illustrate the issues that affect the segment-level evaluation measures used in WMT23, experimentally demonstrating their impact with the help of sentinel metrics; ii) propose solutions to address them; iii) raise concerns regarding the reliability of state-of-the-art MT metrics. We publish the code² to reproduce our work and create novel sentinel metrics.

2 The Meta-evaluation of MT Metrics

Yearly, the WMT Metrics Shared Task organizes a competition among various metrics, including participants' submissions and baselines, with the goal of identifying the metric that most closely aligns with human judgments. Historically, the organizers have employed correlation with human judgment as a meta-evaluation strategy. Recently, significant efforts have been made to refine the meta-evaluation process, encompassing the adoption of new measures, such as those proposed by Kocmi et al. (2021) and Deutsch et al. (2023), and the introduction of the challenge sets sub-task (Freitag et al., 2021b, 2022), among other initiatives. In this section, we provide an overview of WMT's official meta-evaluation setting.

First, multiple MT systems are employed to translate source segments found in one or more test datasets.³ Consequently, test datasets contain several translations of the same source sentence. Second, a manual evaluation campaign is carried out to assess the quality of all translations. Ultimately, metrics' capabilities are assessed based on their alignment with human judgments assigned to translations in the form of scalar scores. Such alignment is typically estimated using correlation and accuracy measures. Specifically, metrics are evaluated at two granularity levels:

- at the segment level, metrics assign a score to every translation of each source segment, and they are ranked according to their ability to discern which translation is superior;
- at the system level, metrics assign a score to each MT system,⁴ and they are ranked according to their ability to discern which system performs better.

At both granularity levels, metrics can be evaluated using several statistical methods, such as the Kendall τ and Pearson ρ correlation coefficients,

²omitted.link ³A segment typically refers to a single sentence, but can also include multiple sentences. For instance, at WMT23, the meta-evaluation of translations from English to German was conducted at the paragraph level rather than at the sentence level.

⁴Typically, the score of a system is calculated as the mean of the scores given to its translations.

which have traditionally been applied at the segment and system levels, respectively. After collecting results from various statistics, a final ranking of metrics is derived by aggregating them. For example, at WMT23, the final ranking was computed from the following three statistics:

- 1. System-level pairwise ranking accuracy (Kocmi et al., 2021), which evaluates metrics based on their ability to rank systems in the same order as human judgments.
- System- and segment-level Pearson correlation, which measures the degree to which metric scores and human scores are linearly correlated.
- 3. Segment-level pairwise ranking accuracy with tie calibration (Deutsch et al., 2023), which evaluates metrics based on their ability to rank segments in the same order as human judgments, or their ability to correctly predict ties.

In this work, we identify two critical issues related to the second and third statistics, and provide the following recommendations to address them:

- Translations should be grouped by their source segment before calculating a metric's segment-level correlation with human judgments (Section 3).
- Tie calibration should not be conducted on the test set (Section 4).

In the following two sections, we provide an overview of the aforementioned statistics, illustrate their flaws, and demonstrate their impact by leveraging our sentinel metrics.

3 To Group or Not to Group?

At early editions of the WMT Metrics Shared Task (Macháček and Bojar, 2013, 2014; Stanojević et al., 2015; Bojar et al., 2016), human gold scores were collected in the form of Relative Rankings (RR). Specifically, the annotators were tasked to rank up to 5 translations of the same source sentence, produced by different MT systems. From each ranking, up to 10 pairwise comparisons were extracted. Despite metrics assessments being scalar scores – which theoretically enable the comparison of all pairs of translated segments – correlation was measured only on those pairs of translations for which RRs were available. Therefore, only translations of the same source sentence were compared. Later on, at more recent editions of WMT, new techniques for human evaluation were adopted: first Direct Assessments (Graham et al., 2013, DA) - where annotators rate individual translations on a scale from 1 to 100 – then Multidimensional Quality Metrics (Lommel et al., 2014, MOM) – where annotators tag the spans of a translation that contain errors, specifying their category and severity. With both the new annotation schemas, each translated segment is assigned a scalar quality score independently of the other segments,⁵ which made it possible to compare all translations, not only those of the same source sentence. This new possibility raised some doubts regarding which is the best way to compute the correlation between metrics and human assessments. Indeed, although both human and metrics assessments can now be represented by matrices - with segments on the x-axis and systems on the y-axis – correlations are applied to vectors, not matrices. Therefore, it was necessary to decide whether to compute the correlation on the flattened matrices - No Grouping - or to first compute the correlations of the rows or columns of the matrices. i.e., group translations based on either their source segment - Segment Grouping - or the system that produced them - System Grouping - respectively, and return the average correlation.

At the WMT21 Metrics Shared Task, Freitag et al. (2021b) chose the *No Grouping* strategy, arguing that the other options would provide only a partial view of the overall picture. At WMT22, all three grouping strategies were used (Freitag et al., 2022), and later at WMT23, Freitag et al. (2023) chose *No Grouping* again. Although *No Grouping* is the only strategy that assesses the MT metrics' ability to discern between higher and lower quality translations in absolute terms, irrespective of the source segment or translation system, we show that both *No Grouping* and *System Grouping* may introduce unfairness and favor trained metrics over the rest.

3.1 The Relation Between Spurious Correlations and Grouping Strategies

Most neural-based metrics are trained with a regression objective to approximate human judgments. They are expected to infer by pattern-matching the relation between human judgments and various

⁵In MQM, a final score is obtained by applying a specific weighting to each combination of the detected spans' category and severity.

phenomena, such as omissions, additions, or other translation errors. However, this mechanism might inadvertently lead to the detection of patterns that are not in a causal relation with the concept of translation quality but are instead spurious correlations, e.g., the number of named entities in the source segment, among others. Arguably, the meta-evaluation should not reward metrics for basing their assessments on spurious correlations between the features of the source, translation, or reference, and the human judgments. However, our intuition is that No Grouping and System Grouping strategies might be doing so by allowing the comparison of translations of different sources. To simplify, consider a metric that unfairly penalizes a translation solely if it contains many named entities. Using No Grouping or System Grouping, such a metric might have a non-negative correlation with human judgments if, on average, translating sentences containing many named entities is more difficult than translating other sentences, because MT systems would be making more mistakes translating them. Therefore, exploiting such a pattern might be beneficial even though it is not causally related to the quality of a translation. In contrast, using Segment Grouping such a pattern would be useless, as translations of the same source sentence should contain the same amount of named entities. More in general, we expect Segment Grouping to lessen the impact of most spurious correlations derived from features shared by a source sentence and its translations.

To assess the extent of this issue in the present evaluation framework, we incorporate three sentinel metrics in the evaluation and re-compute the metrics' rankings using all grouping strategies. Crucially, we find that the impact of spurious correlations when *No Grouping* and *System Grouping* strategies are employed is substantial – favoring trained metrics over the rest⁶ – and is significantly reduced with *Segment Grouping*.

3.2 The Sentinel Metrics

This section describes the three sentinel metrics used to measure the impact of the aforementioned issue on the meta-evaluation process:

1. SENTINEL_{CAND}, which assesses the quality of a translation without comparing it to its source

or any reference.

- 2. SENTINEL_{SRC}, which predicts the quality of a translation solely based on its source.
- 3. SENTINEL_{REF}, which predicts the quality of a translation solely based on its reference.

Having no information regarding the translation to evaluate, SENTINEL_{SRC} and SENTINEL_{REF} can only learn spurious correlations between the features of the source and reference sentences, respectively, and the human judgments. SENTINEL_{CAND}, instead, is a metric with partial information. Indeed, it is possible to evaluate a translation's fluency or grammatical correctness without comparing it with a reference, but not its adequacy. Nonetheless, we expect SENTINEL_{CAND} to base its assessments also on spurious correlations.

3.3 Experimental Setup

Sentinel metrics employ XLM-RoBERTa large (Conneau et al., 2020) as their backbone model, with a multi-layer fully-connected neural network on top of the [CLS] token, which is used to output predictions in the form of scalar scores. Such sentinel metrics are trained by minimizing the Mean Squared Error (MSE) between their predicted quality scores and human judgments. Our dataset comprises a selection of data from WMT spanning 2017 to 2022, incorporating human judgments represented through Direct Assessments (DA) and Multidimensional Quality Metrics (MQM) scores. Inspired by the approach of Rei et al. (2022a), we train sentinel metrics for a single epoch using DA from 2017 to 2020, and then fine-tune them for a further epoch using MQM data. Further details regarding the training process are reported in Appendix B.

3.4 Results

In Table 2, we report the ranking derived from the segment-level Pearson correlation of the primary submissions to the Metrics Shared Task of WMT23, with the inclusion of sentinel metrics, in the language direction $ZH \rightarrow EN$, and with all three grouping strategies. We report in Appendix C the rankings alongside the correlation values for all the other official translation directions of the Metrics Shared Task of WMT23, i.e., $ZH \rightarrow EN$, $EN \rightarrow DE$ and $HE \rightarrow EN$. As can be seen, SENTINEL_{SRC} ranks fourth and third when the grouping strategies are *No Grouping* and *System*

⁶Indeed, overlap-based metrics such as BLEU (Papineni et al., 2002) and chrF (Popović, 2015), or LLM-based metrics such as GEMBA-MQM (Kocmi and Federmann, 2023), were not trained to mimic human assessments and should not be able to leverage spurious correlations.

Grouping, respectively, surpassing strong baselines like COMET or BLEURT-20, and even state-of-theart metrics like GEMBA-MQM. The only metrics that are not surpassed are strong regressors such as XCOMET-Ensemble (Guerreiro et al., 2023) and MetricX-23 (Juraska et al., 2023), which might have learned at least the same spurious correlations leveraged by the sentinel metrics $(\S3.4.1)$. Conversely, when grouping by segment, SENTINEL_{SRC} and SENTINEL_{REF} are correctly positioned at the bottom of the ranking,⁷ and SENTINEL_{CAND} ranks 11th, compared to 3rd and 2nd with No Grouping and System Grouping. A notable difference between the grouping strategies is the positioning of GEMBA-MQM, which is ranked 7th and 9th with No Grouping and System Grouping, respectively, and becomes first with Segment Grouping. We hypothesize that this is due to GEMBA-MQM being based on GPT-4, which has not been explicitly fine-tuned on human assessments and most likely does not leverage spurious correlations such those that described in Section 3.1. Interestingly, with grouping strategies other than Segment Grouping, GEMBA-MQM is surpassed by all the sentinel metrics.

SENTINEL_{CAND} is the only sentinel metric that does not rank at the very bottom with Segment Grouping, outperforming prismSrc (Thompson and Post, 2020) and embed_llama (Dreano et al., 2023), and positioning within the same cluster of statistical significance as BLEU. This suggests that focusing solely on the candidate translation – specifically, its fluency and grammatical correctness - may be sufficient to exceed the performance of some less effective metrics, at least in terms of Pearson correlation with human judgments. Furthermore, we highlight that our results may provide an answer to the open question left at WMT23 regarding the inconsistency of segment-level and system-level correlations for prismSrc. Freitag et al. (2023) noticed that, despite displaying a moderate correlation at the segment level, prismSrc was showing negative correlation values at the system level. As can be seen from Table 2, prismSrc ranks 15th out of 24 with No Grouping but 13th out of 14 with Segment Grouping (i.e., it is in the second to last significance cluster, close to the sentinel metrics). This result is consistent with prismSrc's negative correlation at the system level.

	Grouping				
Metric	No	Seg	Sys		
XCOMET-Ensemble	1	2	1		
MetricX-23-QE*	1	4	1		
XCOMET-QE-Ensemble*	1	3	1		
MetricX-23	2	3	2		
SENTINEL _{CAND} *	3	11	2		
SENTINEL _{SRC} *	4	14	3		
sescoreX	4	7	5		
MaTESe	5	6	6		
SENTINEL _{REF}	5	14	4		
mbr-metricx-qe*	6	1	7		
cometoid22-wmt22*	6	4	6		
GEMBA-MQM*	7	1	9		
Calibri-COMET22-QE*	7	5	8		
CometKiwi*	7	3	9		
KG-BERTScore*	8	4	10		
COMET	9	4	12		
Calibri-COMET22	9	7	11		
docWMT22CometKiwiDA*	10	6	13		
BLEURT-20	10	4	13		
MS-COMET-QE-22*	11	7	14		
docWMT22CometDA	12	6	15		
<u>YiSi-1</u>	13	6	16		
BERTscore	14	7	17		
prismSrc*	15	13	16		
prismRef	16	6	18		
embed_llama	17	12	18		
mre-score-labse-regular	18	8	19		
BLEU	19	11	20		
XLsim	19	10	21		
f200spBLEU	20	10	21		
MEE4	20	9	21		
chrF	21	8	22		
tokengram_F	22	8	23		
Random-sysname*	23	14	23		
eBLEU	24	10	24		

Table 2: Rankings obtained from the segment-level Pearson correlation for the primary submissions to the WMT 2023 Metrics Shared Task, with sentinel metrics. The language direction is $ZH \rightarrow EN$. Ranks represent clusters of statistical significance. Additional information can be found in Appendix C.

In Appendix C, we also report the rankings and correlations obtained using the Kendall τ correlation coefficient for each grouping strategy, to show that our findings are independent of the correlation measure, at least among those typically employed at WMT, i.e., Pearson and Kendall τ .

3.4.1 Are MT metrics learning spurious correlations?

We hypothesize that some of the trained metrics may be basing their assessments on the same spurious correlations leveraged by the sentinel metrics. To delve deeper into this, we measure their segment-level Pearson correlation with the sentinel metrics using *No Grouping*. Surprisingly, XCOMET-Ensemble, XCOMET-QE-Ensemble,

⁷This had to be expected, given that both these metrics return the same assessment for all translations of the same source segment.

MetricX-23, and MetricX-23-QE, which are the only metrics that surpass the sentinels in Table 2, display a high correlation with all sentinel metrics. Interestingly, their correlation with SENTINEL_{SRC} is 0.750, 0.736, 0.690, and 0.712 (Figure 2), respectively, while their correlation with human judgment is 0.650, 0.647, 0.625, and 0.647, respectively (Table 5). We recognize that these metrics share many similarities with our sentinels, as both are neural transformer-based systems and both were trained with the same regression-based objective, using largely the same data. This similarity likely contributes to the high correlation values observed. However, with access limited to only the source segment, SENTINEL_{SRC} relies exclusively on spurious correlations to conduct the evaluation. For this reason, we argue that these results raise concerns about the reliability of state-of-the-art MT metrics, which may be learning to exploit spurious correlations to minimize the Mean Squared Error with human judgments during training. We leave the investigation of this phenomenon to future work and direct readers to Appendix D, where we report the pairwise correlation between most of the considered metrics, for further details.

4 The Evaluation of Ties

In this Section, we focus on the third statistic among those described in Section 2, i.e., the segment-level pairwise ranking accuracy with tie calibration, dubbed acc_{eq} by Deutsch et al. (2023). Before WMT23, the organizers of the Metrics Shared Task used to employ the Kendall τ coefficient – which is a statistic used to estimate the rank-based agreement between two sets of measurements (Kendall, 1945) - to measure the correlation between metrics and human judgments at the segment level. Deutsch et al. (2023) pointed out that the Kendall τ coefficient does not account for metrics correctly predicting ties,⁸ and introduced acceq to address this issue. Unfortunately, our analysis indicates that acceq inadvertently compromises evaluation fairness to accommodate ties, ultimately biasing the results in favor of continuous metrics⁹

over discrete ones.

4.1 The Kendall τ

In this section, we define the Kendall τ coefficient as employed by the organizers of the Metrics Shared Task of WMT21 and WMT22.¹⁰ Let m, h be the vectors of metric and human assessments, respectively. Concordant pairs are the pairs of metric assessments that have been ranked in the same order by humans; discordant pairs are those ranked in a different order. We define C and Das the number of concordant and discordant pairs, respectively. We also define T_h as the number of pairs only tied in the gold scores, T_m as the number of pairs only tied in the metric scores, and T_{hm} as the number of pairs tied both in gold and metric scores, i.e., the number of correctly predicted ties. The Kendall τ correlation coefficient is defined as follows (Kendall, 1945):

$$\tau = \frac{C - D}{\sqrt{(C + D + T_h)(C + D + T_m)}}.$$
 (1)

In Appendix E, we provide a numerical example of the computation of Kendall τ from the vectors m and h.

4.2 The acc_{eq}

As noted by Deutsch et al. (2023), Kendall τ penalizes the prediction of ties, but never rewards them, as T_m and T_h are in the denominator, and T_{hm} is not used. This issue was not prominent in the earliest editions of the Metrics Shared Task, where ties in human scores were disregarded, and older metrics rarely produced ties. Currently, instead, it is essential to consider the prediction of ties, especially since human MQM annotations contain a lot of ties,¹¹ and some recently-proposed metrics are designed to output evaluation assessments that resemble MQM (Perrella et al., 2022; Kocmi and Federmann, 2023). For this reason, Deutsch et al. (2023) proposed a measure that mimics the τ coefficient in the way it is computed, but also accounts for correctly predicting ties:

$$\operatorname{acc}_{\operatorname{eq}} = \frac{C + T_{hm}}{C + D + T_h + T_m + T_{hm}}.$$
 (2)

⁸Given a pair of translations whose quality has been assessed by human annotators, the pair is tied if both translations were assigned with the same score.

⁹By continuous, we refer to those metrics whose assessments can take on any value within a given range, as opposed to discrete metrics, which can take on a limited set of values. Metrics from the COMET family such as COMET, XCOMET-Ensemble, and CometKiwi (Rei et al., 2022b) are continuous, whereas GEMBA-MQM (Kocmi and Federmann, 2023) and

MaTESe (Perrella et al., 2022) are examples of discrete metrics.

¹⁰This is τ_b in Deutsch et al. (2023).

¹¹This is also related to the increasing quality of automatic translation, as perfect translations are assigned the same maximum score.

This measure, as it stands, would unfairly disadvantage most neural metrics, and, in general, continuous metrics. Indeed, it is extremely infrequent for them to assign the same score to two different translations, meaning that they never predict ties. This places them at a disadvantage when human ties are considered in the evaluation. To address this issue, Deutsch et al. (2023) propose the tie calibration algorithm. In the following section, we briefly illustrate such an algorithm and explain why it should not be conducted on the same test set used for the meta-evaluation.

4.3 Tie Calibration

The tie calibration algorithm determines, for each metric, a threshold ϵ such that, given two metric assessments m_1 and m_2 , they are tied if $|m_1$ $m_2 \leq \epsilon$. Deutsch et al. (2023) propose to select the ϵ that maximizes $\operatorname{acc}_{\operatorname{eq}}$ on the same test set used for the metrics meta-evaluation, enabling metrics to output the number of tied scores that best fits the distribution of human ties in the considered test set. This distribution is not stable across test sets (Table 11), and Deutsch et al. (2023) show that ϵ values are not stable either. Nonetheless, they argue that this would not impact the fairness of the evaluation. Unfortunately, our analysis shows that this is not the case. Specifically, despite all metrics' ϵ values being selected on the same test data, we demonstrate that continuous metrics are more flexible to best fit the underlying distribution of human ties, compared to discrete ones, leading to higher acc_{eq} values.

4.4 Two New Sentinel Metrics

To demonstrate the impact of this phenomenon, we introduce two additional sentinel metrics, i.e., SENTINEL_{GEMBA} and SENTINEL_{MATESE}. GEMBA-MQM (Kocmi and Federmann, 2023) and MaTESe (Perrella et al., 2022) are MT metrics that output discrete scores in the form of MQM quality assessments and participated in WMT23. SENTINEL_{GEMBA} and SENTINEL_{MATESE} are perturbed versions of GEMBA-MQM and MaTESe, respectively, obtained by adding Gaussian noise $-\mathcal{N}(0, 0.01)$ – to their predictions. Our objective is to make their output continuous in the neighborhood of discrete values while preventing two different discrete assessments from inverting their ordering. That is, if two GEMBA-MQM's assessments m_1, m_2 are such that $m_1 > m_2$, this relation is preserved by SENTINEL_{GEMBA}. In this way, we create two sentinel metrics that try to partially fill the gap between discrete and continuous metrics, to use them in our analysis. Nonetheless, we wish to remark that this solution is sub-optimal, and is not comparable to metrics that are continuous by design. Indeed, Gaussian noise randomizes the ordering of all SENTINEL_{GEMBA} and SENTINEL_{MATESE}'s assessments that are in the neighborhood of discrete values.

To demonstrate that SENTINELGEMBA and SENTINEL_{MATESE} can better fit the distribution of human ties compared to their discrete counterparts, we modify such distribution in the test data. Specifically, we repeatedly sub-sample the test data, such that for each pair of tied human assessments we remove that pair from the test data with a certain probability p_t , and do the same for non-tied pairs, which are removed with probability p_n . We extract 13 samples by assigning various values to p_t and p_n and report the chosen values in Table 12 in the Appendix. As a consequence, each pair (p_t, p_n) represents a different sub-sample of test data, with a different percentage of tied human pairs. Then, for each metric, we select the best ϵ and compute acc_{eq} on each of these samples.

4.5 Results

In Figure 1a, we present the acc_{eq} results for a subset of continuous metrics, together with GEMBA-MQM, MaTESe, SENTINEL_{GEMBA}, and SENTINEL_{MATESE}. We discuss our results on the WMT23 ZH \rightarrow EN test set, and report results concerning the other language directions, i.e., EN \rightarrow DE and HE \rightarrow EN, in Appendix F. At first glance, it is evident that discrete metrics exhibit a distinct acc_{eq} pattern compared to continuous and sentinel metrics. Notably, at lower percentages of tied human pairs, SENTINEL_{GEMBA} and SENTINEL_{MATESE} significantly outperform GEMBA-MQM and MaTESe.¹² This discrepancy arises because the tie calibration algorithm selects

¹²It is important to highlight that the range of human tie percentages explored in our analysis is similar to that found in the WMT test sets. Indeed, as shown in Table 11, such percentages range from a minimum of 15.14% to a maximum of 53.35%, observed in the WMT22 EN \rightarrow DE test set.



Figure 1: $\operatorname{acc_{eq}}(a)$ and optimal ϵ (b) of the considered metrics for varying percentages of human ties in the test dataset (0.24 is the percentage of human ties in the entire dataset, obtained when p_t and p_n are both 0). ϵ values have been scaled using min-max scaling. Specifically, for each metric, the minimum ϵ is the optimal ϵ at 0% of human ties, and the maximum is the optimal ϵ at 100%. The language direction is ZH \rightarrow EN. Results concerning all language directions can be found in Appendix F. For each percentage of human ties, we use 5 different seeds to sub-sample the test data. Therefore, the shown $\operatorname{acc_{eq}}$ and ϵ , for each metric and percentage of ties, are averaged across 5 different runs. Best seen in color.

very small ϵ values, close to 0 for every metric, allowing the number of ties predicted by continuous metrics to potentially drop to 0. Conversely, metrics that yield discrete scores inherently produce a certain number of ties, placing them at a disadvantage, and thus ranking conceptually identical metrics like SENTINELGEMBA and GEMBA-MQM at significantly different positions. Interestingly, in the hypothetical scenario in which there are no tied human pairs in the dataset, SENTINELGEMBA would rank second (despite lots of its assessments having a random ordering), whereas GEMBA-MQM would be second to last. At increasing percentages of gold ties, instead, the acc_{eq} values obtained by SENTINEL_{GEMBA} and SENTINEL_{MATESE} converge to those of their discrete counterparts. However, this is a limitation of these sentinels' design and does not imply that the evaluation is fair at higher percentages of human ties. We delve deeper into this matter in Figure 1b, which shows how the optimal ϵ changes at varying percentages of human ties. As can be seen from the figure, continuous metrics' ϵ is dynamically adjusted with heightened sensitivity, contrary to what happens for discrete and sentinel metrics. Specifically, their ϵ is exactly 0 until 39% of human ties. Additionally, for MaTESe, it

remains constant between 44% and 56%, and between 61% and 68%, and the same happens for GEMBA-MQM between 47% and 51% and between 56% and 68%. In contrast, the values change for all the other metrics in the same intervals.

5 Conclusion

In this work, we identified two issues with the current meta-evaluation of Machine Translation metrics, as conducted at the Metrics Shared Task of the Conference on Machine Translation. We proposed a suite of sentinel metrics designed to highlight these issues and demonstrate their impact on the metrics rankings, revealing that the current metaevaluation process tends to favor certain metric categories. Specifically, the None Grouping and System Grouping strategies prefer trained metrics over overlap- and LLM-based ones, and the algorithm of tie optimization, if conducted on the same test set used for the meta-evaluation, favors continuous metrics over discrete ones. Furthermore, we observed a notably high correlation between sentinel metrics and state-of-the-art metrics, raising concerns about their reliability and suggesting that their assessments might be based on spurious correlations present in the training data.

Limitations

We recognize that the *Segment Grouping* approach does not evaluate the ability of metrics to distinguish between higher and lower quality translations in absolute terms, that is, independently of their source sentence. This aspect should indeed play a role in the meta-evaluation process. However, our analysis suggests that the rankings derived from the *No Grouping* and *System Grouping* methods favor certain metric categories, and potentially reward metrics for leveraging spurious correlations. Nonetheless, we believe that there is a need to develop fairer methods to fill this gap in the metaevaluation.

Regarding the optimization of ϵ in acc_{eq}, our analysis does not definitively specify the optimal selection process for ϵ values. Although we demonstrated that continuous metrics are favored by selecting the optimal ϵ on the same test set used for the meta-evaluation, this does not necessarily mean that using a held-out dataset ensures a fair metaevaluation. The distribution of human ties in the held-out dataset could either advantage or disadvantage continuous metrics, due to their greater adaptability in fitting such distribution, compared to discrete metrics. In this respect, Appendix F presents the acceq score of MT metrics, calculated after applying the tie calibration algorithm to a sub-sample of the test set and then computing acceq across the entire test set, rather than just the sub-sample as discussed in Section 4.5. We observe that continuous metrics are at a disadvantage when the proportion of ties in the sample used to estimate ϵ significantly deviates from the proportion in the entire test set. In this context, the flexibility of continuous metrics to adapt to the underlying distribution of human ties acts as a drawback, rather than a benefit. In general, we believe that a promising approach might involve estimating statistically significant score deltas for continuous metrics, treating as tied all assessments within these deltas (akin to the work of Kocmi et al. (2024) regarding system-level assessments). This approach would also enhance the interpretability of MT metrics.

Acknowledgements

References

Ondřej Bojar, Yvette Graham, Amir Kamran, and Miloš Stanojević. 2016. Results of the WMT16 metrics shared task. In Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Pa*pers*, pages 199–231, Berlin, Germany. Association for Computational Linguistics.

- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale.
- Daniel Deutsch, Rotem Dror, and Dan Roth. 2021. A statistical analysis of summarization evaluation metrics using resampling methods. *Transactions of the Association for Computational Linguistics*, 9:1132–1146.
- Daniel Deutsch, George Foster, and Markus Freitag. 2023. Ties matter: Meta-evaluating modern metrics with pairwise accuracy and tie calibration. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12914– 12929, Singapore. Association for Computational Linguistics.
- Sören Dreano, Derek Molloy, and Noel Murphy. 2023. Embed_Llama: Using LLM embeddings for the metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 738–745, Singapore. Association for Computational Linguistics.
- Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021a. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.
- Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.
- Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference* on Machine Translation (WMT), pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021b. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, pages 733–774, Online. Association for Computational Linguistics.

- Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. Continuous measurement scales in human evaluation of machine translation. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 33–41, Sofia, Bulgaria. Association for Computational Linguistics.
- 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, 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 *Proceedings* of the Eighth Conference on Machine Translation, pages 756–767, Singapore. Association for Computational Linguistics.
- M. G. Kendall. 1945. The treatment of ties in ranking problems. *Biometrika*, 33:239–51.
- 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.
- Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes. 2021. To ship or not to ship: An extensive evaluation of automatic metrics for machine translation. In *Proceedings of the Sixth Conference on Machine Translation*, pages 478–494, Online. Association for Computational Linguistics.
- Tom Kocmi, Vilém Zouhar, Christian Federmann, and Matt Post. 2024. Navigating the metrics maze: Reconciling score magnitudes and accuracies.
- Philipp Koehn and Christof Monz. 2006. Manual and automatic evaluation of machine translation between European languages. In *Proceedings on the Work-shop on Statistical Machine Translation*, pages 102– 121, New York City. Association for Computational Linguistics.
- Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. 2021. On the variance of the adaptive learning rate and beyond.
- Arle. Language Technology Lab) Lommel, Hans. Language Technology Lab) Uszkoreit, and Aljoscha. Language Technology Lab) Burchardt. 2014. Multidimensional quality metrics (mqm) : a framework for declaring and describing translation quality metrics. *Tradumàtica*, (12):455–463.
- Matouš Macháček and Ondřej Bojar. 2013. Results of the WMT13 metrics shared task. In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 45–51, Sofia, Bulgaria. Association for Computational Linguistics.

- Matouš Macháček and Ondřej Bojar. 2014. Results of the WMT14 metrics shared task. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 293–301, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. 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.
- 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, Lisbon, Portugal. Association for 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 *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, 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.
- 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 *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.
- Miloš Stanojević, Amir Kamran, Philipp Koehn, and Ondřej Bojar. 2015. Results of the WMT15 metrics shared task. In *Proceedings of the Tenth Workshop*

on Statistical Machine Translation, pages 256–273, Lisbon, Portugal. Association for Computational Linguistics.

- Brian Thompson and Matt Post. 2020. 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.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc.

A Overall ranking

In Table 3, we report the official segment-level ranking of WMT23 Metrics Shared Task, including sentinel metrics.

Training the sentinel metrics B

The input for the sentinel metrics consists of either source text (SENTINEL_{SRC}), candidate translation (SENTINEL_{CAND}), or reference translation (SENTINEL_{REF}). This text is first tokenized, and then passed to the XLM-RoBERTa large model, serving as the backbone for feature extraction. From this, the embedding of the [CLS] token is extracted to use it as a comprehensive representation of the input. Then, the [CLS] token embedding is fed into a multi-layer fully-connected neural network, which outputs the scalar quality score. More formally, considering t as the input text for a sentinel metric:

$$\begin{split} \boldsymbol{e}_t &= \operatorname{XLMR}\left(t\right) \\ \boldsymbol{h}_t^{(1)} &= \operatorname{Dropout}\left(\operatorname{Tanh}\left(W_h^{(1)}\boldsymbol{e}_t + \boldsymbol{b}_h^{(1)}\right)\right) \\ \boldsymbol{h}_t^{(2)} &= \operatorname{Dropout}\left(\operatorname{Tanh}\left(W_h^{(2)}\boldsymbol{h}_t^{(1)} + \boldsymbol{b}_h^{(2)}\right)\right) \\ s_t &= W_o\boldsymbol{h}_t^{(2)} + \boldsymbol{b}_o \end{split}$$

Where:

- e_t is the embedding of the [CLS] token from the input text t, obtained through the XLM-RoBERTa large model.
- $m{h}_t^{(i)}$ represents the output of the i^{th} layer, with each layer consisting of a linear transformation (using weight matrix $W_h^{(i)}$ and bias vector $b_{h}^{(i)}$) followed by an activation function and regularization.

- W_o and \boldsymbol{b}_o are the weight matrix and bias vector, respectively, for the output layer.
- s_t is the scalar quality score.

Both the initial training phase using DA data and the subsequent fine-tuning phase with MQM scores employ the same set of hyperparameters, detailed in Table 4.

Grouping Strategies С

In Tables 5, 6, 7, we report the complete set of rankings and Pearson correlations, at the segment level, of the primary submission to the WMT23 Metrics Shared Task, with sentinel metrics. Sentinel metrics are consistently ranked lower with Segment Grouping. However, this grouping strategy requires the estimation of multiple Pearson correlation coefficients - one for each group of translations – which are ultimately averaged. As a consequence, the number of clusters of statistical significance is reduced.

In Tables 8, 9, 10, we report the complete set of rankings and Kendall τ correlation coefficients, at the segment level, of the primary submissions to the WMT23 Metrics Shared Task, with sentinel metrics. With Kendall τ as well, sentinel metrics have worse performances when Segment Grouping is employed. We wish to remark that one should not focus on the absolute values of the ranks but on their value relative to that of the other metrics. Indeed, as already mentioned, the number of clusters of statistical significance is reduced with Segment Grouping. For instance, in Table 9, SENTINELCAND is ranked 5th out of 19 with No Grouping, and 4th out of 11 with Segment Grouping. While the absolute value of the rank is higher, in terms of correlation it has moved from the 8th to the 17th position.

Metrics pairwise correlations D

In Figures 2, 3, 4, we report the pairwise correlation between a part of the primary submissions and baselines of WMT23, with the inclusion of sentinel metrics. We use Pearson correlation coefficient, with No Grouping. As can be seen, state-of-the-art regression-based metrics display a notably high correlation with sentinels. Specifically, the highest correlations are reported for XCOMET-Ensemble and MetricX-23, and their reference-less counterparts. Moderate correlation is also reported between sentinels and baseline metrics such as CometKiwi,

			${ m EN} ightarrow { m DE}$			${ m HE} ightarrow { m EN}$				${ m ZH} ightarrow { m EN}$				
Metric	Avg	g. corr	Р	earson	a	cc _{eq}	Pe	arson	a	cc _{eq}	P	earson	a	cc _{eq}
XCOMET-Ensemble	1	0.697	1	0.695	1	0.604	1	0.556	1	0.586	1	0.650	1	0.543
MetricX-23	2	0.682	4	0.585	1	0.603	1	0.548	2	0.577	2	0.625	3	0.531
XCOMET-QE-Ensemble*	3	0.681	2	0.679	3	0.588	3	0.498	4	0.554	1	0.647	3	0.533
MetricX-23-QE*	4	0.681	3	0.626	2	0.596	2	0.520	3	0.564	1	0.647	4	0.527
mbr-metricx-qe*	5	0.652	4	0.571	3	0.584	5	0.411	4	0.553	6	0.489	2	0.537
GEMBA-MQM*	6	0.639	6	0.502	5	0.572	5	0.401	3	0.564	7	0.449	5	0.522
MaTESe	7	0.636	5	0.554	9	0.528	4	0.459	5	0.550	5	0.511	12	0.479
CometKiwi*	8	0.632	7	0.475	5	0.569	7	0.387	6	0.544	7	0.442	4	0.525
sescoreX	9	0.628	6	0.519	6	0.563	7	0.385	16	0.484	4	0.536	9	0.499
SENTINEL _{CAND} *	10	0.626	5	0.561	6	0.562	10	0.339	16	0.483	3	0.580	14	0.473
cometoid22-wmt22*	11	0.625	8	0.441	4	0.578	9	0.365	12	0.515	6	0.479	7	0.515
KG-BERTScore*	12	0.624	8	0.451	7	0.556	8	0.382	7	0.537	8	0.430	6	0.516
COMET	13	0.622	9	0.432	4	0.574	5	0.401	8	0.532	9	0.396	7	0.514
BLEURT-20	14	0.622	7	0.484	5	0.572	8	0.382	11	0.519	10	0.378	6	0.518
Calibri-COMET22-QE*	15	0.603	9	0.441	12	0.483	6	0.395	13	0.506	7	0.443	10	0.491
Calibri-COMET22	16	0.603	10	0.413	10	0.522	5	0.401	12	0.515	9	0.396	14	0.474
YiSi-1	17	0.600	12	0.366	8	0.542	6	0.395	8	0.529	12	0.290	8	0.504
docWMT22CometDA	18	0.598	11	0.394	7	0.559	10	0.339	14	0.497	11	0.353	10	0.493
docWMT22CometKiwiDA*	19	0.598	8	0.444	8	0.547	12	0.286	15	0.489	9	0.387	10	0.493
prismRef	20	0.593	6	0.516	10	0.518	11	0.319	9	0.528	14	0.183	8	0.504
MS-COMET-QE-22*	21	0.588	13	0.310	8	0.546	12	0.295	14	0.498	10	0.367	9	0.498
BERTscore	22	0.582	13	0.325	9	0.528	10	0.335	12	0.515	13	0.236	9	0.499
mre-score-labse-regular	23	0.558	18	0.111	9	0.530	8	0.378	10	0.522	16	0.145	12	0.481
XLsim	24	0.544	14	0.239	9	0.527	14	0.233	17	0.480	17	0.111	15	0.464
f200spBLEU	25	0.540	14	0.237	9	0.526	14	0.230	19	0.447	18	0.108	13	0.476
MEE4	26	0.539	17	0.202	9	0.529	13	0.256	20	0.441	18	0.105	12	0.480
tokengram_F	27	0.537	16	0.227	10	0.520	14	0.226	18	0.461	20	0.060	11	0.485
chrF	28	0.537	15	0.232	10	0.519	15	0.221	18	0.460	19	0.063	11	0.485
BLEU	29	0.533	17	0.192	10	0.520	15	0.220	20	0.442	17	0.119	14	0.472
prismSrc*	30	0.530	9	0.425	13	0.426	16	0.140	20	0.441	13	0.223	17	0.421
embed_llama	31	0.529	14	0.250	12	0.483	15	0.215	21	0.430	15	0.161	16	0.447
SENTINEL _{SRC} *	32	0.512	7	0.469	15	0.231	10	0.334	21	0.428	4	0.540	19	0.240
SENTINEL _{REF}	33	0.506	8	0.464	15	0.231	11	0.301	21	0.428	5	0.506	19	0.240
eBLEU	34	0.491	20	-0.011	11	0.512	16	0.131	19	0.445	22	-0.084	14	0.473
Random-sysname*	35	0.463	19	0.064	14	0.409	17	0.041	21	0.428	21	0.018	18	0.381

Table 3: Complete segment-level results for the primary submissions	s to the	WMT	2023 M	etrics S	Shared	Task,	with
sentinel metrics.							

COMET, and BLEURT-20. As expected, instead, close to no correlation is reported for lexical-based metrics such as BLEU and chrF, which are not trained metrics. Similarly, GEMBA-MQM, a stateof-the-art LLM-based metric that has not been finetuned on human assessments, shows low levels of correlation with the sentinel metrics.

Е Kendall τ and acc_{eq} computation

In this section, we provide an example of the computation of Kendall τ and acc_{eq} from two vectors of human and metric scores, i.e., h and m in the following table:

m	0.6	0.5	0.4	0.4
h	5	3	5	5

For example, the pairs of metric assessments are (m_1, m_2) , (m_1, m_3) , (m_1, m_4) , (m_2, m_3) , $(m_2, m_4), (m_3, m_4).$

In Equations 1 and 2, C = 1, since the only concordant pair is (m_1, m_2) . Indeed, $m_1 >$ m_2 and $h_1 > h_2$. D = 2, since the pairs $(m_2, m_3), (m_2, m_4)$ are discordant. $T_m = 0$, since there are no pairs tied only in the metric scores. $T_h = 2$, since the pairs $(h_1, h_3), (h_1, h_4)$ are tied only in the human scores. $T_{hm} = 1$, since the remaining pair, i.e., (m_3, m_4) , is tied in both human and metric scores. In this example, $\tau = -0.258$ and $acc_{eq} = 0.333$.

F Ties

In Table 11, we report the percentage of tied human pairs in several datasets containing human judgments in the form of MQM scores.

For each vector, there are six pairs of assessments.



Figure 2: Pairwise correlation between a part of the primary submissions and baselines of WMT23, and sentinel metrics. Correlation is Pearson with *No Grouping*, and the language direction is $ZH \rightarrow EN$.



Figure 3: Pairwise correlation between a part of the primary submissions and baselines of WMT23, and sentinel metrics. Correlation is Pearson with *No Grouping*, and the language direction is $EN \rightarrow DE$.



Figure 4: Pairwise correlation between a part of the primary submissions and baselines of WMT23, and sentinel metrics. Correlation is Pearson with *No Grouping*, and the language direction is $HE \rightarrow EN$.

Hyperparameter	Value
Optimizer	RAdam (Liu et al., 2021)
Learning Rate	1e-6
Number of Epochs	1
Batch Size	8
Accumulation Steps	2
Dropout	0.1
Dimension of $oldsymbol{h}_t^{(1)}$	512
Dimension of $oldsymbol{h}_t^{(2)}$	128

Table 4: Hyperparameters used for training and finetuning the sentinel metrics.

In Tables 12, 13, 14, we report the values of p_t and p_n used to sub-sample the ZH \rightarrow EN, EN \rightarrow DE, and HE \rightarrow EN test sets, respectively, to conduct the experiment illustrated in Section 4.4. We also report the corresponding percentage of human ties and total number of pairs, for each sample.

In Figures 5a, 6a, 7a, we report the acc_{eq} and optimal ϵ for each of the considered metrics, in all three language directions considered at WMT 2023.

In Figure 8, we report the acc_{eq} values of the considered metrics on the entire test set. ϵ values have been estimated from sub-samples of the test data, each sub-sample having a different percentage of human ties (represented by the numbers on the x-axis).

Motric	No		Sa	amont	System		
Metric		INU	36	gment			
XCOMET-Ensemble	1	0.650	2	0.421	1	0.610	
MetricX-23-QE*	1	0.647	4	0.359	1	0.610	
XCOMET-QE-Ensemble*	1	0.647	3	0.380	1	0.612	
MetricX-23	2	0.625	3	0.373	2	0.580	
SENTINEL _{CAND} *	3	0.580	11	0.201	2	0.578	
SENTINEL _{SRC} *	4	0.540	14	0.000	3	0.561	
sescoreX	4	0.536	7	0.295	5	0.505	
MaTESe	5	0.511	6	0.325	6	0.441	
SENTINEL _{REF}	5	0.506	14	0.000	4	0.525	
mbr-metricx-qe*	6	0.489	1	0.436	7	0.431	
cometoid22-wmt22*	6	0.479	4	0.357	6	0.446	
GEMBA-MQM*	7	0.449	1	0.434	9	0.378	
Calibri-COMET22-QE*	7	0.443	5	0.355	8	0.411	
<u>CometKiwi</u> *	7	0.442	3	0.388	9	0.388	
KG-BERTScore*	8	0.430	4	0.369	10	0.374	
COMET	9	0.396	4	0.364	12	0.345	
Calibri-COMET22	9	0.396	7	0.311	11	0.360	
docWMT22CometKiwiDA*	10	0.387	6	0.340	13	0.320	
BLEURT-20	10	0.378	4	0.371	13	0.330	
MS-COMET-QE-22*	11	0.367	7	0.306	14	0.313	
docWMT22CometDA	12	0.353	6	0.327	15	0.291	
<u>YiSi-1</u>	13	0.290	6	0.329	16	0.237	
BERTscore	14	0.236	7	0.309	17	0.186	
prismSrc*	15	0.223	13	0.078	16	0.243	
prismRef	16	0.183	6	0.332	18	0.135	
embed_llama	17	0.161	12	0.138	18	0.139	
mre-score-labse-regular	18	0.145	8	0.251	19	0.123	
BLEU	19	0.119	11	0.208	20	0.093	
XLsim	19	0.111	10	0.218	21	0.069	
f200spBLEU	20	0.108	10	0.220	21	0.077	
MEE4	20	0.105	9	0.236	21	0.070	
chrF	21	0.063	8	0.263	22	0.020	
tokengram_F	22	0.060	8	0.262	23	0.015	
Random-sysname*	23	0.018	14	0.019	23	0.002	
eBLEU	24	-0.084	10	0.219	24	-0.115	

Table 5: Segment-level Pearson correlation for the primary submissions to the WMT23 Metrics Shared Task, with sentinel metrics. The language direction is $ZH \rightarrow EN$. Starred metrics are reference-free, underlined metrics are baselines, and highlighted metrics are sentinels. Ranks represent clusters of statistical significance and are computed following Freitag et al. (2023), which leverage the PERM-BOTH hypothesis test introduced by Deutsch et al. (2021).

Metric		No	Se	gment	System		
XCOMET-Ensemble	1	0.695	1	0.538	1	0.676	
XCOMET-QE-Ensemble*	2	0.679	2	0.507	2	0.658	
MetricX-23-OE*	3	0.626	2	0.511	3	0.564	
MetricX-23	4	0.585	2	0.507	4	0.547	
mbr-metricx-qe*	4	0.571	1	0.543	3	0.551	
SENTINEL _{CAND} *	5	0.561	6	0.396	5	0.522	
MaTESe	5	0.554	8	0.330	4	0.526	
sescoreX	6	0.519	3	0.459	6	0.502	
prismRef	6	0.516	7	0.349	4	0.528	
GEMBA-MQM*	6	0.502	3	0.482	$\overline{7}$	0.446	
BLEURT-20	7	0.484	2	0.492	$\overline{7}$	0.455	
CometKiwi*	$\overline{7}$	0.475	3	0.463	7	0.451	
SENTINEL _{SRC} *	8	0.469	12	0.000	6	0.502	
SENTINEL _{REF}	8	0.464	12	0.000	6	0.492	
KG-BERTScore*	8	0.451	4	0.456	8	0.421	
docWMT22CometKiwiDA*	9	0.444	5	0.426	9	0.404	
cometoid22-wmt22*	9	0.441	2	0.499	9	0.385	
Calibri-COMET22-QE*	9	0.441	5	0.432	8	0.414	
<u>COMET</u>	9	0.432	2	0.508	10	0.363	
prismSrc*	9	0.425	11	0.102	6	0.487	
Calibri-COMET22	10	0.413	3	0.477	10	0.370	
docWMT22CometDA	11	0.394	3	0.484	11	0.310	
<u>YiSi-1</u>	12	0.366	5	0.404	12	0.284	
<u>BERTscore</u>	13	0.325	7	0.355	13	0.250	
MS-COMET-QE-22*	13	0.310	6	0.400	13	0.241	
embed_llama	14	0.250	10	0.242	14	0.180	
XLsim	14	0.239	6	0.372	16	0.151	
f200spBLEU	14	0.237	7	0.343	14	0.178	
chrF	15	0.232	8	0.336	15	0.157	
tokengram_F	16	0.227	8	0.340	16	0.153	
MEE4	17	0.202	7	0.360	16	0.145	
BLEU	17	0.192	9	0.310	17	0.140	
mre-score-labse-regular	18	0.111	6	0.376	18	0.087	
Random-sysname*	19	0.064	11	0.124	19	-0.015	
eBLEU	20	-0.011	8	0.317	19	-0.030	

Table 6: Segment-level Pearson correlation for the primary submissions to the WMT23 Metrics Shared Task, with sentinel metrics. The language direction is $EN \rightarrow DE$. Starred metrics are reference-free, underlined metrics are baselines, and highlighted metrics are sentinels. Ranks represent clusters of statistical significance and are computed following Freitag et al. (2023), which leverage the PERM-BOTH hypothesis test introduced by Deutsch et al. (2021).

Metric	No		Se	gment	System		
XCOMET-Ensemble	1	0.556	1	0.479	1	0.515	
MetricX-23	1	0.548	2	0.441	1	0.509	
MetricX-23-QE*	2	0.520	5	0.387	2	0.480	
XCOMET-QE-Ensemble*	3	0.498	4	0.397	3	0.458	
MaTESe	4	0.459	5	0.373	4	0.408	
mbr-metricx-qe*	5	0.411	2	0.448	5	0.362	
GEMBA-MQM*	5	0.401	2	0.431	6	0.354	
COMET	5	0.401	3	0.421	5	0.367	
Calibri-COMET22	5	0.401	4	0.397	5	0.371	
<u>YiSi-1</u>	6	0.395	2	0.439	6	0.348	
Calibri-COMET22-QE*	6	0.395	6	0.354	5	0.369	
CometKiwi*	7	0.387	5	0.375	6	0.353	
sescoreX	7	0.385	5	0.370	6	0.352	
KG-BERTScore*	8	0.382	5	0.375	7	0.347	
BLEURT-20	8	0.382	3	0.418	7	0.344	
mre-score-labse-regular	8	0.378	4	0.407	8	0.335	
cometoid22-wmt22*	9	0.365	7	0.309	7	0.346	
docWMT22CometDA	10	0.339	5	0.379	9	0.294	
SENTINEL _{CAND} *	10	0.339	11	0.104	7	0.343	
BERTscore	10	0.335	4	0.412	9	0.293	
SENTINEL _{SRC} *	10	0.334	13	0.000	7	0.336	
prismRef	11	0.319	3	0.428	10	0.276	
SENTINEL _{REF}	11	0.301	13	0.000	9	0.299	
MS-COMET-QE-22*	12	0.295	9	0.252	10	0.274	
docWMT22CometKiwiDA*	12	0.286	7	0.324	11	0.234	
MEE4	13	0.256	8	0.291	11	0.222	
XLsim	14	0.233	7	0.314	12	0.198	
f200spBLEU	14	0.230	8	0.287	12	0.195	
tokengram_F	14	0.226	$\overline{7}$	0.311	13	0.184	
<u>chrF</u>	15	0.221	$\overline{7}$	0.308	14	0.179	
BLEU	15	0.220	9	0.260	13	0.189	
embed_llama	15	0.215	10	0.188	13	0.187	
prismSrc*	16	0.140	11	0.100	15	0.150	
eBLEU	16	0.131	8	0.280	16	0.104	
Random-sysname*	17	0.041	12	0.057	17	0.001	

Table 7: Segment-level Pearson correlation for the primary submissions to the WMT23 Metrics Shared Task, with sentinel metrics. The language direction is $HE \rightarrow EN$. Starred metrics are reference-free, underlined metrics are baselines, and highlighted metrics are sentinels. Ranks represent clusters of statistical significance and are computed following Freitag et al. (2023), which leverage the PERM-BOTH hypothesis test introduced by Deutsch et al. (2021).

Metric		No	Se	Segment		ystem
XCOMET-Ensemble	1	0.473	2	0.299	1	0.456
XCOMET-QE-Ensemble*	2	0.467	3	0.273	2	0.451
MetricX-23-QE*	3	0.461	4	0.252	2	0.448
GEMBA-MQM*	3	0.457	1	0.365	4	0.416
MetricX-23	4	0.449	3	0.269	3	0.434
mbr-metricx-qe*	5	0.427	2	0.301	5	0.403
cometoid22-wmt22*	5	0.423	4	0.252	4	0.408
SENTINEL _{CAND} *	6	0.404	9	0.148	4	0.410
SENTINEL _{SRC} *	7	0.397	14	0.000	4	0.411
CometKiwi*	7	0.391	3	0.263	6	0.368
Calibri-COMET22-QE*	8	0.386	4	0.241	6	0.366
sescoreX	9	0.375	6	0.217	6	0.367
MaTESe	9	0.371	3	0.271	7	0.345
KG-BERTScore*	10	0.361	4	0.248	8	0.337
SENTINEL _{REF}	11	0.340	14	0.000	7	0.353
COMET	11	0.333	4	0.248	9	0.311
MS-COMET-QE-22*	11	0.332	6	0.213	9	0.311
Calibri-COMET22	12	0.330	6	0.217	9	0.310
BLEURT-20	13	0.310	3	0.261	10	0.288
docWMT22CometKiwiDA*	14	0.299	5	0.234	11	0.265
docWMT22CometDA	15	0.276	5	0.231	12	0.248
prismSrc*	16	0.234	12	0.044	12	0.251
YiSi-1	17	0.220	5	0.231	13	0.196
BERTscore	18	0.180	6	0.216	14	0.156
mre-score-labse-regular	18	0.178	7	0.176	14	0.165
prismRef	19	0.165	5	0.232	15	0.140
embed_llama	20	0.109	11	0.096	16	0.093
XLsim	20	0.101	10	0.140	17	0.080
MEE4	21	0.091	8	0.172	18	0.064
<u>BLEU</u>	21	0.085	9	0.154	18	0.062
f200spBLEU	22	0.068	8	0.165	19	0.042
chrF	23	0.045	7	0.187	20	0.017
tokengram_F	24	0.042	$\overline{7}$	0.187	21	0.012
Random-sysname*	25	0.015	13	0.025	22	-0.005
eBLEU	26	-0.041	9	0.156	23	-0.064

Table 8: Segment-level Kendall τ correlation coefficient for the primary submissions to the WMT23 Metrics Shared Task, with sentinel metrics. The language direction is $ZH \rightarrow EN$. Starred metrics are reference-free, underlined metrics are baselines, and highlighted metrics are sentinels. Ranks represent clusters of statistical significance and are computed following Freitag et al. (2023), which leverage the PERM-BOTH hypothesis test introduced by Deutsch et al. (2021).

Metric		No	Se	gment	S	ystem
XCOMET-Ensemble	1	0.546	1	0.380	1	0.530
XCOMET-QE-Ensemble*	2	0.532	2	0.360	2	0.516
MetricX-23-QE*	3	0.509	2	0.357	3	0.487
MetricX-23	3	0.506	2	0.368	3	0.485
sescoreX	4	0.493	3	0.343	4	0.476
mbr-metricx-qe*	4	0.490	1	0.397	4	0.467
GEMBA-MQM*	4	0.482	1	0.399	5	0.449
SENTINEL _{CAND} *	5	0.463	4	0.290	5	0.456
MaTESe	5	0.462	5	0.286	6	0.447
BLEURT-20	6	0.452	2	0.366	7	0.426
SENTINEL _{SRC} *	6	0.443	11	0.000	5	0.462
cometoid22-wmt22*	7	0.422	2	0.362	8	0.398
SENTINEL _{REF}	7	0.418	11	0.000	6	0.437
COMET	7	0.418	2	0.366	9	0.387
Calibri-COMET22	7	0.417	3	0.342	9	0.387
CometKiwi*	8	0.408	3	0.330	9	0.379
Calibri-COMET22-QE*	8	0.406	5	0.279	9	0.379
MS-COMET-QE-22*	9	0.391	5	0.280	10	0.363
KG-BERTScore*	10	0.361	4	0.310	11	0.329
docWMT22CometKiwiDA*	10	0.358	4	0.316	11	0.329
prismRef	11	0.345	6	0.247	11	0.332
docWMT22CometDA	11	0.337	2	0.360	12	0.296
<u>YiSi-1</u>	12	0.280	4	0.297	13	0.250
prismSrc*	12	0.267	10	0.039	12	0.284
BERTscore	13	0.253	5	0.260	14	0.224
MEE4	14	0.225	5	0.271	15	0.190
XLsim	14	0.217	6	0.257	15	0.180
f200spBLEU	15	0.187	6	0.255	16	0.151
chrF	15	0.186	6	0.241	16	0.152
tokengram_F	16	0.183	6	0.245	17	0.149
embed_llama	16	0.182	8	0.163	16	0.150
BLEU	17	0.137	$\overline{7}$	0.231	18	0.103
eBLEU	18	0.096	$\overline{7}$	0.230	19	0.070
mre-score-labse-regular	18	0.084	5	0.269	19	0.066
Random-sysname*	19	0.033	9	0.081	20	-0.018

Table 9: Segment-level Kendall τ correlation coefficient for the primary submissions to the WMT23 Metrics Shared Task, with sentinel metrics. The language direction is EN \rightarrow DE. Starred metrics are reference-free, underlined metrics are baselines, and highlighted metrics are sentinels. Ranks represent clusters of statistical significance and are computed following Freitag et al. (2023), which leverage the PERM-BOTH hypothesis test introduced by Deutsch et al. (2021).

Metric		No	Se	gment	Sy	stem
XCOMET-Ensemble	1	0.415	2	0.323	1	0.395
MetricX-23	2	0.401	3	0.302	2	0.382
GEMBA-MQM*	2	0.399	1	0.369	3	0.367
XCOMET-QE-Ensemble*	3	0.374	5	0.276	3	0.358
MetricX-23-QE*	3	0.370	6	0.251	3	0.355
mbr-metricx-qe*	3	0.366	2	0.316	4	0.339
MaTESe	4	0.361	3	0.302	4	0.341
COMET	5	0.350	3	0.309	5	0.327
Calibri-COMET22	6	0.348	4	0.284	6	0.324
BLEURT-20	6	0.344	4	0.295	6	0.320
sescoreX	6	0.342	4	0.285	6	0.320
CometKiwi*	7	0.338	6	0.238	6	0.323
Calibri-COMET22-QE*	7	0.336	7	0.230	6	0.322
<u>YiSi-1</u>	7	0.333	2	0.325	7	0.303
mre-score-labse-regular	7	0.328	4	0.284	7	0.300
KG-BERTScore*	8	0.322	6	0.242	7	0.304
cometoid22-wmt22*	9	0.310	7	0.216	7	0.301
prismRef	9	0.302	3	0.309	8	0.273
BERTscore	10	0.295	4	0.298	9	0.266
docWMT22CometDA	11	0.278	5	0.270	10	0.249
MS-COMET-QE-22*	12	0.261	9	0.174	10	0.249
SENTINEL _{SRC} *	13	0.243	12	0.000	10	0.247
SENTINEL _{CAND} *	13	0.243	11	0.049	10	0.249
XLsim	13	0.233	7	0.228	11	0.211
MEE4	13	0.231	7	0.221	11	0.202
docWMT22CometKiwiDA*	14	0.227	7	0.229	12	0.192
SENTINEL _{REF}	15	0.210	12	0.000	11	0.214
tokengram_F	15	0.207	7	0.228	13	0.175
<u>chrF</u>	16	0.204	7	0.224	14	0.171
f200spBLEU	17	0.193	7	0.219	15	0.162
BLEU	18	0.184	8	0.205	16	0.157
embed_llama	18	0.174	10	0.147	16	0.151
eBLEU	19	0.166	8	0.209	17	0.141
prismSrc*	19	0.164	11	0.043	14	0.169
Random-sysname*	20	0.027	11	0.033	18	0.002

Table 10: Segment-level Kendall τ correlation coefficient for the primary submissions to the WMT23 Metrics Shared Task, with sentinel metrics. The language direction is HE \rightarrow EN. Starred metrics are reference-free, underlined metrics are baselines, and highlighted metrics are sentinels. Ranks represent clusters of statistical significance and are computed following Freitag et al. (2023), which leverage the PERM-BOTH hypothesis test introduced by Deutsch et al. (2021).



Figure 5: $\operatorname{acc}_{eq}(a)$ and optimal ϵ (b) of the considered metrics for varying percentages of human ties in the test dataset (0.24 is the percentage of human ties in the entire dataset, obtained when p_t and p_n are both 0). ϵ values have been scaled using min-max scaling. Specifically, for each metric, the minimum ϵ is the optimal ϵ at 0% of human ties, and the maximum is the optimal ϵ at 100%. The language direction is ZH \rightarrow EN. For each percentage of human ties, we use 5 different seeds to sub-sample the test data. Therefore, the shown acc_{eq} and ϵ , for each metric and percentage of ties, are averaged across 5 different runs.



Figure 6: $\operatorname{acc}_{eq}(a)$ and optimal ϵ (b) of the considered metrics for varying percentages of human ties in the test dataset (0.23 is the percentage of human ties in the entire dataset, obtained when p_t and p_n are both 0). ϵ values have been scaled using min-max scaling. Specifically, for each metric, the minimum ϵ is the optimal ϵ at 0% of human ties, and the maximum is the optimal ϵ at 100%. The language direction is EN \rightarrow DE. For each percentage of human ties, we use 5 different seeds to sub-sample the test data. Therefore, the shown acc_{eq} and ϵ , for each metric and percentage of ties, are averaged across 5 different runs.



Figure 7: $\operatorname{acc}_{eq}(a)$ and optimal ϵ (b) of the considered metrics for varying percentages of human ties in the test dataset (0.43 is the percentage of human ties in the entire dataset, obtained when p_t and p_n are both 0). ϵ values have been scaled using min-max scaling. Specifically, for each metric, the minimum ϵ is the optimal ϵ at 0% of human ties, and the maximum is the optimal ϵ at 100%. The language direction is HE \rightarrow EN. For each percentage of human ties, we use 5 different seeds to sub-sample the test data. Therefore, the shown acc_{eq} and ϵ , for each metric and percentage of ties, are averaged across 5 different runs.





(a) The language pair is ZH \rightarrow EN. 0.24 is the percentage of human ties in the sub-sample used to estimate ϵ .

(b) The language pair is EN \rightarrow DE. 0.23 is the percentage of human ties in the sub-sample used to estimate $\epsilon.$



(c) The language pair is HE \to EN. 0.43 is the percentage of human ties in the sub-sample used to estimate $\epsilon.$

Figure 8: acc_{eq} of the considered metrics on the entire test set. For each metric, ϵ values are estimated using sub-samples of the test set, with varying percentages of human ties, that are on the x-axis. For each percentage of human ties, we use 5 different seeds to sub-sample the test data. Therefore, the shown acc_{eq} for each metric and percentage of ties are averaged across 5 different runs.

	2020	2021	2022	2023
$EN \rightarrow DE$	15.14	44.62	53.35	23.11
$ZH \rightarrow EN$ $EN \rightarrow RU$	17.01	$\frac{30.31}{53.24}$	$\begin{array}{c} 41.55\\ 44.42\end{array}$	24.03
$\rm HE \rightarrow EN$	_	_	_	42.84

Table 11: Percentage of tied pairs in the MQM data released over different years at the Metrics Shared Task (or by Freitag et al. (2021a), for 2020), and regarding different translation directions.

p_t	1.00	0.65	0.30	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
p_n	0.00	0.00	0.00	0.00	0.20	0.40	0.50	0.60	0.65	0.70	0.75	0.80	0.85
%	0	10	18	24	28	35	39	44	47	51	56	61	68
#	93890	104304	114664	123585	104888	85969	76522	67237	62624	57948	53110	48491	43730

Table 12: p_t is the probability of removing a tied human pair, and p_n is that of removing a non-tied human pair. The considered test set is WMT23 ZH \rightarrow EN. Each column, i.e., each pair (p_t, p_n) , represents a sub-sample of the test set, in which tied and non-tied pairs have been removed with such probabilities. The third row contains the percentage of tied human pairs over all pairs, as a result of the sub-sampling. The last row contains the total number of pairs remaining in the test set after the sub-sampling.

p_t	1.00	0.65	0.30	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
p_n	0.00	0.00	0.00	0.00	0.20	0.40	0.50	0.60	0.65	0.7	0.75	0.80	0.85
%	0	10	17	23	27	33	38	43	46	50	54	60	67
#	23343	25803	28236	30360	25694	21021	18689	16353	15184	14014	12899	11698	10493

Table 13: p_t is the probability of removing a tied human pair, and p_n is that of removing a non-tied human pair. The considered test set is WMT23 EN \rightarrow DE. Each column, i.e., each pair (p_t, p_n) , represents a sub-sample of the test set, in which tied and non-tied pairs have been removed with such probabilities. The third row contains the percentage of tied human pairs over all pairs, as a result of the sub-sampling. The last row contains the total number of pairs remaining in the test set after the sub-sampling.

p_t	1.0	0.90	0.80	0.65	0.50	0.35	0.20	0.00	0.00	0.00	0.00	0.00	0.00
p_n	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.20	0.40	0.55	0.65	0.75
%	0	7	13	21	27	33	38	43	48	56	62	68	75
#	36561	39254	42038	46202	50272	54435	58516	63960	56679	49315	43918	40145	36530

Table 14: p_t is the probability of removing a tied human pair, and p_n is that of removing a non-tied human pair. The considered test set is WMT23 HE \rightarrow EN. Each column, i.e., each pair (p_t, p_n) , represents a sub-sample of the test set, in which tied and non-tied pairs have been removed with such probabilities. The third row contains the percentage of tied human pairs over all pairs, as a result of the sub-sampling. The last row contains the total number of pairs remaining in the test set after the sub-sampling.