

# Contextual Metric Meta-Evaluation by Measuring Local Metric Accuracy

Anonymous ACL submission

## Abstract

Meta-evaluation of automatic evaluation metrics—assessing evaluation metrics themselves—is crucial for accurately benchmarking natural language processing systems and has implications for scientific inquiry, production model development, and policy enforcement. While existing approaches to metric meta-evaluation focus on general statements about the absolute and relative quality of metrics across arbitrary system outputs, in practice, metrics are applied in highly contextual settings, often measuring the performance for a highly constrained set of system outputs. For example, we may only be interested in evaluating a specific model or class of models. We introduce a method for contextual metric meta-evaluation by comparing the *local metric accuracy* of evaluation metrics. Across translation, speech recognition, and ranking tasks, we demonstrate that the local metric accuracies vary both in absolute value and relative effectiveness as we shift across evaluation contexts.

## 1 Introduction

Meta-evaluation of automatic evaluation metrics—assessing evaluation metrics themselves—is crucial for accurately benchmarking natural language processing systems (Zhou et al., 2022). Because metrics are central to scientific inquiry, production model development, and policy enforcement (Kocmi et al., 2021), there is a constant need for new approaches to evaluating system outputs (Novikova et al., 2017).

Although current methods for metric meta-evaluation commonly take a *global* perspective, reporting the performance of a metric across arbitrary system outputs, coming from any system (Stanojević et al., 2015; Przybocki et al., 2009), practical metric meta-evaluation is highly contextual, measuring the performance for a highly constrained set of system outputs. For example, we may only be

metric	context			global
	X	Y	Z	
A	0.9	0.9	0.3	0.7
B	0.7	0.7	0.7	0.7
C	0.3	0.3	0.9	0.5

Table 1: Contextual metric meta-evaluation. When comparing metrics A, B, and C, traditional meta-evaluation focuses on global accuracy across arbitrary inputs. Local metric accuracy can vary by evaluation contexts X, Y, and Z.

interested in evaluating a specific model or class of models. From a model development perspective, we may be interested in a metric that is sensitive to model outputs coming from partially trained models at the beginning of the development cycle (when the outputs are far from the target distribution or close to random); such a metric may struggle to differentiate between outputs from fully trained or more effective models. This is highly reflective of the results found by Fomicheva and Specia (2019), who show that metric performance varies significantly across different levels of translation quality. Thus, using the same metric throughout the development process may lead to biased or incomplete evaluations and possibly pruning earlier models, which may have a better performance when fully trained.

To illustrate the difference between global and contextual metric meta-evaluation, we constructed a toy meta-evaluation for three metrics across three contexts (Table 1). The values in the table represent the accuracy of three metrics (A, B, and C) under three different contexts (X, Y, and Z) as well as the global accuracy across the different contexts. By looking at the average, we might think that A and B are equally accurate. However, when inspecting accuracy within individual contexts, we can see that selecting the most appropriate metric is far less straightforward. For example, if we want

071	a metric that best generalizes across different contexts, we want to pick B over A even though their global accuracies are equal. However, if we want to specifically measure outputs in context Z, then we would want to pick C as it is especially sensitive to system outputs in that context, despite it having the lowest global accuracy. This suggests that comparing a metric on its global accuracy may deviate from local accuracy, which may be more relevant in a contextual setting.	121
072		122
073		123
074		124
075		125
076		126
077		127
078		128
079		129
080		130
081	To improve contextual metric meta-evaluation, we propose analyzing metrics across different evaluation contexts and measuring their <i>local metric accuracies</i> . By evaluating metrics across three different machine learning tasks—machine translation, automated speech recognition, ranking—we show that the metric accuracy, which measures the ability of a metric to accurately assign the true preference between a pair of system decisions, changes as the context of the output space changes. We also show that the metric accuracy changes both in absolute value and relative ordering across the different contexts. In contrast with existing work on metric meta-evaluation relies heavily on costly and time-consuming explicit human feedback (Fabri et al., 2021; Liu et al., 2016), our method uses output perturbations (Sai et al., 2021; He et al., 2023) to obtain the true ordering between a pair of system outputs without the need of human supervision. Overall, we show that measuring local metric accuracies is a straightforward methodology to provide a more contextual understanding of evaluation metrics which complements existing global metric meta-evaluation methods.	131
082		132
083		133
084		134
085		135
086		136
087		137
088		138
089		139
090		140
091		141
092		142
093		143
094		144
095		145
096		146
097		147
098		148
099		149
100		150
101		151
102		152
103		153
104		154
105	<b>2 Related work</b>	
106	Our work connects to the broader literature on meta-evaluation, which has been approached in various ways, highlighting the complexity and the necessity of this task. For example, the Workshop of Statistical Machine Translation (WMT) has focused on evaluating the utility of metrics in machine translation since 2008, where participants submit automated metrics for validation against human feedback (Callison-Burch et al., 2008). However, human feedback can be subjective and susceptible to social biases (Sun et al., 2022). Noting these limitations, Xiao et al. (2023) propose a theory-driven meta-evaluation framework rooted in measurement theory for NLG metrics. Their work highlights issues in human evaluation including a	155
107		156
108		157
109		158
110		159
111		160
112		161
113		162
114		163
115		164
116		165
117		166
118		167
119		168
120		169
		170
	lack of validation, standardization, and consistency.	
	Our use of output perturbation is inspired by prior work in testing metric robustness. Chen and Eger (2023) proposed a preference-based adversarial attack framework using targeted perturbations to evaluate the robustness of NLI-based and BERT-based metrics, finding that NLI-based metrics are more robust in summarization but not in machine translation. Sai et al. (2021) extends perturbation-based robustness testing by creating templates targeting specific criteria such as jumbled word order to test fluency. Chen et al. (2019) assessed QA metrics by converting multiple-choice datasets into free-response formats, highlighting the need for BERT-based metrics. Additionally, Valcarce et al. (2018) evaluated the robustness of ranking metrics against incompleteness by introducing sparsity to system outputs. Our paper adopts similar perturbation techniques to assess the preference-based evaluation capability of different metrics, eliminating the need for costly and time-intensive explicit human feedback.	
	Although metric meta-evaluation is often done on a global level, previous work indicates that the reliability of a metric changes from the system-level to the decision-level (Reiter and Belz, 2009; Stent et al., 2005). Though some research has investigated metric performance for different contexts based on output sources (i.e., models) or output qualities (Mathur et al., 2020; Novikova et al., 2017), our work addresses the lack of a systematic review of contextual meta-evaluation and how to conduct it.	
	<b>3 Local accuracy</b>	
	To formalize local metric accuracy, we introduce the following notation. Let $\mathcal{X}$ be the set of all possible system inputs (e.g., for MT, all possible strings from the source language) and $\mathcal{Y}$ the set of all possible system outputs (e.g., for MT, all possible strings from the target language). We define $X \subset \mathcal{X}$ to be the subset of system inputs observed in a specific context (e.g., for MT, a sample of source sentences from a specific university). Similarly, $Y_x \subset \mathcal{Y}$ is the subset of system decisions for $x \in X$ observed for X in a specific context (e.g., for MT, a set of translations generated by a set of candidate systems). In addition, we have access to a perturbation function that, with high probability, degrades the utility of a decision $y$ (e.g., dropping a random word from a translated input). Let $Q_x$	

171 be the set of pairs decisions conditioned on an input  
 172  $x$  and their corresponding degraded version:  
 173  $Q_x = \{\langle y, y' \rangle\}_{y \in Y_x}$ .

174 An evaluation metric  $\mu : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  gener-  
 175 ates a scalar number reflecting the performance  
 176 according to some system property that we want to  
 177 measure (e.g. correctness of translation). Each met-  
 178 ric is an approximation of  $\mu^*$ , the ideal evaluation  
 179 metric (i.e., the true utility of an output). Given  
 180 any two pairs of system outputs,  $\mu^*$  will always  
 181 be able to determine the true ordering of the two  
 182 outputs. In cases where we intentionally perturb  
 183  $y$  to obtain  $y'$ , we know that  $\mu^*(x, y) > \mu^*(x, y')$ .  
 184 Under the assumption that  $\mu$  approximates  $\mu^*$ , we  
 185 want to compute how often  $\mu(x, y) > \mu(x, y')$ . As  
 186 suggested by Kocmi et al. (2021), we focus on  
 187 the ability of  $\mu$  to reproduce the ordering of deci-  
 188 sions, rather than the magnitude of the difference  
 189 between  $\mu(x, y)$  and  $\mu(x, y')$ . From this, we define  
 190 the pointwise local metric accuracy, conditioned  
 191 on an input  $x$  as,

$$192 \text{Acc}_\mu(Q_x) = \frac{1}{|Q_x|} \sum_{\langle y, y' \rangle \in Q_x} \mathbb{1} [\mu(x, y) > \mu(x, y')] \quad (1)$$

193 This measures the ability of a metric to reproduce  
 194 the true ordering of perturbations for a specific  
 195 input  $x$ . We define the local metric accuracy across  
 196 all contexts as,

$$197 \text{Acc}_\mu(Q) = \frac{1}{|X|} \sum_{x \in X} \text{Acc}_\mu(Q_x) \quad (2)$$

198 where  $Q = \cup_{x \in X} Q_x$ . This measures the local  
 199 metric accuracy across a sample of system inputs,  
 200 as we may have in a standard evaluation set.

201 We are interested in testing two hypotheses with  
 202 respect to local metric accuracy.

203 H1: The absolute local metric accuracy,  $\text{Acc}_\mu(Q)$ ,  
 204 of a metric  $\mu$  changes as the context changes.

205 Evidence supporting this hypothesis suggests that  
 206 existing evaluation methods focusing on global  
 207 metric accuracy obscures how metric accuracy  
 208 varies across different contexts.

209 H2: The ordering of a set of metrics by local metric  
 210 accuracy changes as context changes.

211 In other words, the total ordering of all metrics  
 212 by local metric accuracy within a context changes  
 213 as the context  $Q$  changes. Evidence supporting  
 214 this hypothesis suggests that choosing an appropri-  
 215 ate metric to benchmark compare system outputs  
 216 largely depends on the context.

## 4 Methods and Materials 217

### 4.1 Tasks, dataset, and metrics 218

219 We performed our evaluation on three different  
 220 tasks: Machine Translation (MT), Automated  
 221 Speech Recognition (ASR), and Ranking. Table  
 222 2 details the dataset and metrics that we used in  
 223 our experiments. For each task, we used readily  
 224 available system outputs to improve reproducibil-  
 225 ity. For each metric, we employed their respective  
 226 official implementations or, when unavailable, the  
 227 most widely used implementation with default pa-  
 228 rameters. For any neural metric computation, we  
 229 used a NVIDIA RTX A6000 GPU. For BLEU, we  
 230 used nltk’s sentencebleu implementation. We  
 231 also used nltk’s implementation for METEOR. For  
 232 ROUGE<sup>1</sup> and BERTSCORE<sup>2</sup>, we have used the im-  
 233 plementation released by their respective authors.  
 234 For BLEURT<sup>3</sup>, COMET<sup>4</sup>, CHRf<sup>5</sup> and UNITE<sup>6</sup>, we  
 235 have used their official implementations via the  
 236 evaluate library on HuggingFace. We used the  
 237 jiwer<sup>7</sup> Python package to compute the ASR met-  
 238 rics. We used the trec\_eval<sup>8</sup> to calculate the rank-  
 239 ing metrics.

240 We adopt the category with the highest num-  
 241 ber of contexts for each task. The abundance of  
 242 contexts allowed us to identify trends in metric be-  
 243 havior across a broader range of items and helped  
 244 us identify supporting evidence for or against our  
 245 hypotheses.

### 4.2 Perturbation techniques 246

247 To test our hypotheses, we applied a perturbation  
 248 function that degrades the utility of a system output  
 249  $y$  and its corresponding degraded version  $y'$ . Thus,  
 250 we know that the quality of  $y'$  under a specific task  
 251 is worse than  $y$  with a high probability. For the  
 252 system outputs belonging to the machine transla-  
 253 tion and automated speech recognition tasks, we  
 254 perturbed  $y$  by removing 20% of the words in the  
 255 outputs, rounded to the nearest integer.

256 Our perturbation technique is a simplification of  
 257 He et al. (2023), who synthesize a range of per-  
 258 turbations that closely mimic human or machine

<sup>1</sup><https://github.com/google-research/google-research/tree/master/rouge>

<sup>2</sup>[https://github.com/Tiiiger/bert\\_score](https://github.com/Tiiiger/bert_score)

<sup>3</sup><https://github.com/google-research/bleurtreadme>

<sup>4</sup><https://unbabel.github.io/COMET/html/models.html>

<sup>5</sup><https://github.com/mjpost/sacreBLEUchr-f-chr-f>

<sup>6</sup><https://huggingface.co/Unbabel/unite-mup>

<sup>7</sup><https://github.com/jitsi/jiwer>

<sup>8</sup>[https://github.com/usnistgov/trec\\_eval](https://github.com/usnistgov/trec_eval)

Task	Dataset	Metrics
MT	Over 150,000 system outputs and reference translation from 62 different MT systems submitted to the WMT metrics task from year 2023 (Freitag et al., 2023) for the source-target language pairs English-Russian (en-ru), English-German (en-de), Chinese-English (zh-en). The subsets that are available are YEAR, DOMAIN, and SYSTEM.	BLEU, ROUGE-1, ROUGE-2, ROUGE-L, METEOR, BERTSCOREP, BERTSCORER, BERTSCOREF1, COMET, BLEURT, CHRf, UNITESRC, UNITEREF, UNITEUNIFIED
ASR	Over 33,000 system outputs from six different ASR models on ESPnet (Watanabe et al., 2018) on the LibriSpeech 100 dataset (Panayotov et al., 2015). The subsets that are available are SYSTEM, SPEAKER ID, GENDER, and QUALITY.	Word Error Rate (WER), Match Error Rate (MER), Word Information Lost (WIL), Word Information Preserved (WIP), Character Error Rate (CER)
Ranking	Ranked list of top-100 items retrieved by 21 recommender algorithms provided by Valcarce et al. (2018) on the MovieLens1M dataset (Harper and Konstan, 2015) submitted to TREC (Buckley and Voorhees, 2004). We were able to segment the outputs by ALGORITHM.	Mean Average Precision (MAP), Precision@ $R$ , where $R$ is the number of relevant documents (RPREC), Reciprocal Rank (RECIP_RANK), Interpolated Precision at Recall Level $X$ (for $X = \{0.0, 0.1, 0.2, 0.3, 0.4\}$ ) (IPREC_AT_RECALL_X), Precision@ $K$ (P_K), Recall@ $K$ (RECALL_K), nDCG@ $K$ (NDCG_CUT_K) (where $K = 5, 10, 15, 20, 30$ )

Table 2: Datasets and metrics used for different tasks

errors. However, we want added simplicity and generalizability to languages other than English, so we refrained from doing perturbations that are semantically informed, such as removing articles and prepositions, verb lemmatization, or negation.

For the system outputs belonging to the ranking task, we perturbed  $y$  by shuffling the rankings within the top-100 items for each user-system pair.

### 4.3 Hypothesis Testing

In order to test H1, we plotted the metric accuracies  $ACC_\mu(Q)$  for each task across different contexts within the selected context category  $Q$  as a line graph, such that we can visualize how the metric’s capability of differentiating between  $y$  and  $y'$  changes as the context changes by observing the slopes and overlaps between the lines. To further investigate the association between the context  $Q$  and the metric accuracy  $ACC_\mu(Q)$ , we used the  $\chi^2$  test of independence of variables (Pearson, 1900) in a contingency table (Pearson, 1904). We will compare the resulting p-values to the significance level of  $\alpha = 0.05$  to understand whether the changes in metric accuracy  $ACC_\mu(Q)$  across the different contexts  $Q$  are statistically significant.

To test H2, we computed the Kendall’s  $\tau$  (Kendall, 1938) between a two rankings of metrics according to local metric accuracy under two contexts. This helps us quantify how the total ordering of the metrics changes as the context changes. In order to emphasize the metric selection task, we adopt version of Kendall’s  $\tau$  that weighs changes

at the higher in the ranking more than those lower in the ranking (Shieh, 1998). Specifically, we use a hyperbolic weighing that maps each rank  $r$  to weight  $\frac{1}{r+1}$ .

## 5 Results

### 5.1 Machine Translation

We visualize the metric accuracies for the machine translation metrics under the different SYSTEM contexts, as shown in Figure 1a. We observe that the line for each metric changes as we change the context, as indicated by the varying slopes of the lines. The results of the  $\chi^2$  test indicate that the difference in the metric accuracies across the different context is statistically significant, supporting H1 for MT.

Figure 1a also contains intersections between lines corresponding to different metrics, indicating that there is a change in the relative position of each metric in the different contexts and hence signifying a change in the total ordering of metrics by local accuracy across contexts, supporting H2. This is further strengthened in Figure 1b, where the  $\tau$  values show that the correspondence between the pairs of metric accuracy rankings varies considerably for each pair of SYSTEMS. If the  $\tau$  values were consistently close to 1, there would not be support for H2. Instead, we find that  $\tau$  values cluster according to similarity of context.

## 5.2 Automated Speech Recognition

For ASR, we report the local metric accuracy under the different SPEAKER IDs which come from different dataset QUALITY contexts (CLEAN/RAW). We plot the local metric accuracy for contexts associated with different contexts shown in Figure 2a. We observe that the lines corresponding to each metric are not straight, which indicates that the absolute local accuracy for each metric changes with context, supporting H1. Our  $\chi^2$  test results confirm that the difference in the local metric accuracies across the different contexts is statistically significant, providing evidence supporting H1 for ASR.

Interestingly, we do not observe the same overlap between the lines corresponding to the different metrics as we did for MT. This, along with the consistently high values of  $\tau$  in Figure 2b, indicates that there is not evidence supporting H2 for ASR.

## 5.3 Ranking

Plotting the metric accuracies for the ranking metrics for the different ALGORITHMS in Figure 3a, we can first observe that none of the lines corresponding to the different metrics are straight lines, which supports H1. The large fluctuations within each line suggest that the changes in the absolute local accuracies for each metric are rather significant. The  $\chi^2$  test results shows a statistically significant change, providing evidence in support of H1.

Furthermore, we can see that overlaps exist between the different lines corresponding to the different metrics, similar to the observation we made in the MT case. This indicates that the total ordering of metric accuracies changes as the context changes, supporting H2. The  $\tau$  results (Figure 3b) show clustering by algorithm, as with MT.

# 6 Discussion

## 6.1 H1: Absolute Local Accuracies

The results in Section 5 generally provide evidence supporting H1, as our experiments consistently show that the local metric accuracy changes as the context changes.

We can observe that the local metric accuracy for a context is related to the average quality of outputs in that context. For example, in the ranking setting, the most effective system according to MAP is SLIM, which is also the context whose perturbed outputs are easiest to distinguish. Conversely, perturbed outputs in the random ranker are

more difficult for all metrics to distinguish. This is because our perturbation method catastrophically degrades good outputs and bad outputs are already poor and difficult to make demonstrably worse. We will return to this in Section 6.3.

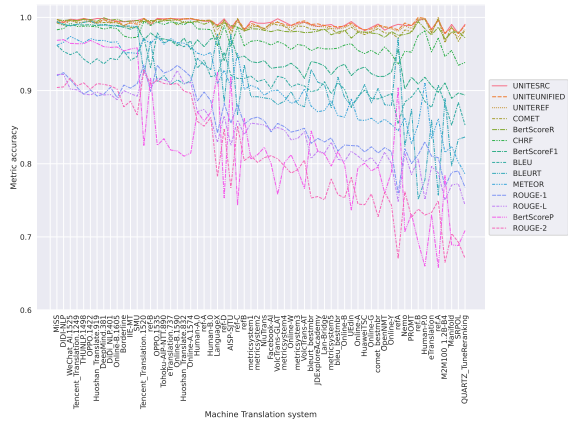
These results suggest that evaluators may be interested in the stability of local metric accuracy when selecting a metric. A metric that is more stable with respect to local metric accuracy is more predictable when deployed under a new context and the probability of selecting the wrong system is consistent. This is especially important if we consider a new evaluation context where poor local metric accuracy puts users—or a vulnerable subgroup of users—at risk. Work in robust machine learning provides existing methods for designing metrics stable across context changes (Yuan et al., 2024).

In addition to stability, we can organize metrics according to systematic behavior in local metric accuracy. For example, in Figure 1a, more complex embedding-based and model-based evaluation metrics generally perform better than the simpler lexical-based metrics (Zhang et al., 2019; Freitag et al., 2022). More complex metrics cluster on the top of the figure while simpler metrics occupy the bottom regions; any overlap occurs mostly within a specific category. Such analysis allows evaluators to understand the empirical relationships between metric ensembles. Although picking the best metric might involve selecting a metric occupying the top of the figure, there may be contexts in which local metric accuracies are close enough to allow flexibility in selecting metrics with lower local metric accuracy.

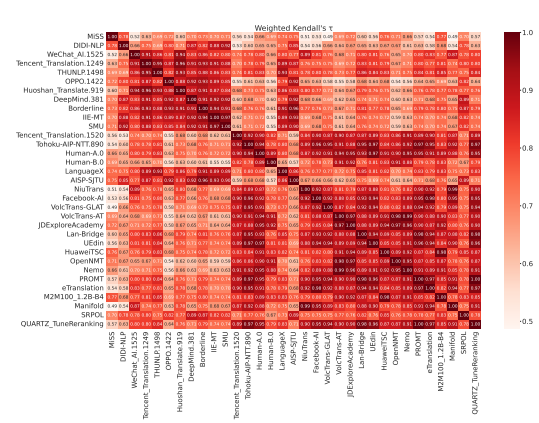
More generally, we can consider multi-objective metric development. For example, since embedding- and model-based methods are more time-intensive and computationally costly compared to the lexical-based methods, adopting simpler and cheaper metrics when local metric accuracies are comparable (e.g., early in model development) would result in cost savings and faster iteration. Beyond cost and local metric accuracy, one can imagine local versions of metric interpretability, metric engineering overhead, metric optimizability, and other criteria when conducting for contextual meta-evaluation.

## 6.2 H2: Relative Local Accuracies

Although the observations in Section 5 errs toward accepting H2, the evidence from our experi-

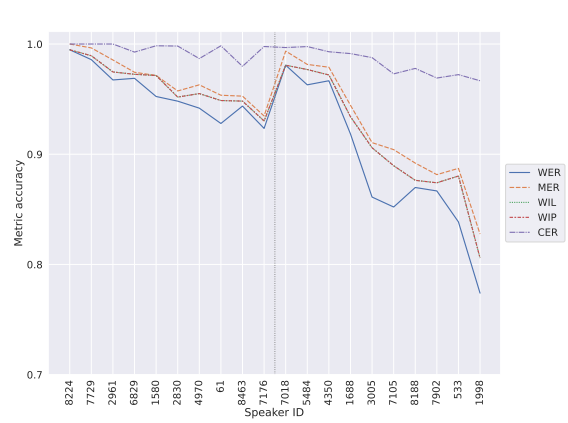


(a) Local metric accuracy across the different contexts

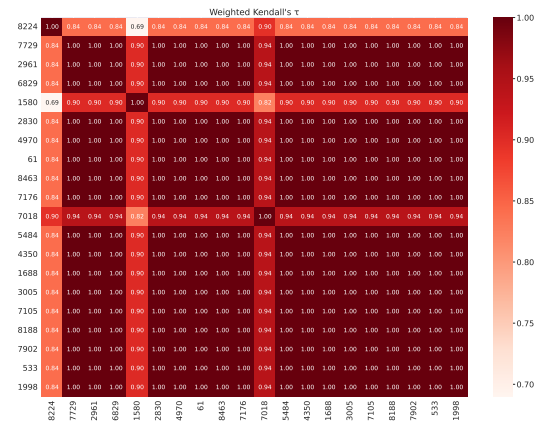


(b) Weighted Kendall's  $\tau$  of metrics ordered by local accuracy between different contexts

Figure 1: Machine Translation. Contexts provided by translation system.

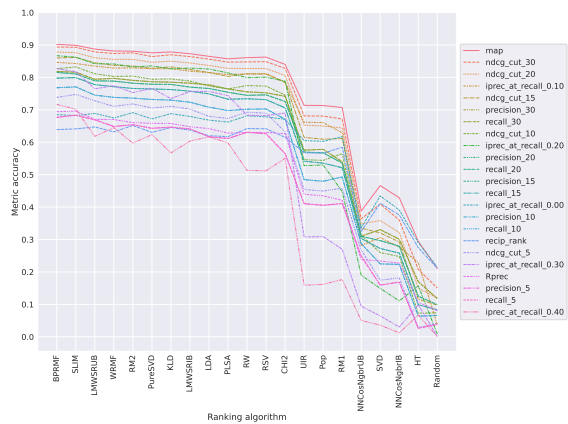


(a) Local metric accuracy across the different contexts

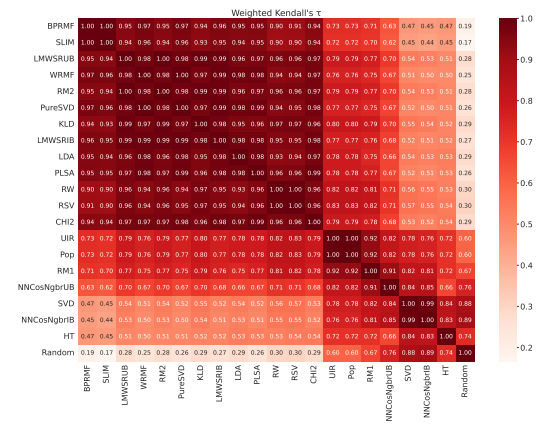


(b) Weighted Kendall's  $\tau$  of metrics ordered by local accuracy between different contexts

Figure 2: Automatic Speech Recognition. Local metric accuracy across different Speaker IDs. (a) Speaker IDs to the left of the grey line come from the QUALITY=CLEAN LibriSpeech-100 dataset, while the Speaker IDs to the right of the grey line come from the QUALITY=OTHER LibriSpeech-100 dataset.



(a) Local metric accuracy across the different contexts



(b) Weighted Kendall's  $\tau$  of metrics ordered by local accuracy between different contexts

Figure 3: Ranking. Metric accuracy for Ranking metrics across the different systems.

ments on the ASR task in Section 5.2 suggest that it strongly depends on the nature of the task and the metrics. Evaluating ASR is relatively straightforward, where the ambiguity of correct answers is low, unlike MT or ranking where two outputs (e.g., translations or permutations) can be equally good (Wieting et al., 2019). Hence, the metrics are commonly used to benchmark ASR systems only slightly vary in the construct they are trying to measure, and they are all operationalized following similar statistical methods.

Figures 1b and 3b indicate that there are groups of contexts where the relative reliability of metrics is similar. When contexts can be structured according to metric accuracy, ones can adopt a fixed evaluation metric. This has practical implications in terms of engineering and development overhead or, in the case of model-based metrics, model development cost. Predicting the similarity in local metric accuracy ordering (i.e., the cells in Figures 1b and 3b) is an important task because it allows evaluators to confidently adopt an evaluation metric without conducting contextual meta-evaluation. Predictive features include any metadata we have about the contexts. For example, in ranking, Valcarce et al. (2018) categorize ALGORITHMS into different families of techniques: matrix factorization (SVD, PURESVD, BPRMF, WRMF), neighborhood-based (CHI2, KLD, RSV, ROCCHIO’S WEIGHTS).

### 6.3 Methodology

Although our results demonstrate that local metric accuracy analysis can provide insight into metric behavior, there are several opportunities for improving the methodology. First, our perturbations, while reliable in generating output degradations, may result in outputs that are easily detected by metrics, especially for highly effective systems. Moreover, perturbed outputs may be sufficiently different as to be unlikely to occur in a specific context. For example, if we are evaluating in the context of highly effective MT systems, a translation with a missing word is very unlikely by any highly effective MT system, even though we know it is lower quality. In order to address this, developing perturbation methods that reliably degrade performance and are likely to occur within a context will be important for future local metric accuracy development. This is related to synthesizing hard negative examples in the contrastive learning literature (Kalantidis et al., 2020). Alternatively, we

can consider non-perturbation data, perhaps from human annotators, although this compromises the cost-effectiveness of output perturbation.

In order to help with clarity, we focused on contexts that were interpretable, which contexts are relevant depends on the broader model evaluation environment. Focusing on models, as we did for MT and ranking, emphasizes contexts that reflect iterative model development and refinement within a narrow set of constraints (i.e., the particular model being evaluated). If we are benchmarking a diverse set of systems, we are interested in comparing a broader set of possible outputs than those from a single system. In cases where we are designing a metric agnostic to a particular context, we may be interested in robust performance across arbitrary contexts. While this is similar to global analysis, a more rigorous and formal approach to context selection, such as found in the distributionally robust machine learning literature (Duchi et al., 2018), may be more appropriate.

## 7 Conclusion

We introduce the notion of local metric accuracy and demonstrate how to use it to conduct contextual metric meta-evaluation. Our results show that both the absolute and relative local accuracy of a metric varies as we vary context, though this depends on the nature of the task. Based on our results, we believe that by moving beyond global metric meta-evaluation, we can achieve a more accurate understanding of metric performance, which in turn increases the reliability of the evaluations and provide actionable insights for improving NLP systems.

## 8 Limitations

As mentioned in Section 4.2, our experiments adopted relatively simple perturbation methods in order to cover a wide range of tasks and guarantee degradation. In future work, we plan to explore more task- and language-specific methods developed in the NLP community (Sai et al., 2021; Chen and Eger, 2023; He et al., 2023).

We also compute local accuracy by uniformly weighting all output-perturbation pairs. In reality, different outputs have different probabilities of occurring in a specific context. These probabilities should be incorporated into the accuracy calculation to provide a more reliable estimate of local metric accuracy. Estimating the distribution over

518	outputs for a specific context itself is a difficult		
519	research question which we plan on addressing in		
520	future work.		
	<b>References</b>		
521			
522	Chris Buckley and Ellen M Voorhees. 2004. Retrieval		
523	evaluation with incomplete information. In <i>Proceed-</i>		
524	<i>ings of the 27th annual international ACM SIGIR</i>		
525	<i>conference on Research and development in informa-</i>		
526	<i>tion retrieval</i> , pages 25–32.		
527	Chris Callison-Burch, Philipp Koehn, Christof Monz,		
528	Josh Schroeder, and Cameron Shaw Fordyce. 2008.		
529	Proceedings of the third workshop on statistical ma-		
530	chine translation. In <i>Proceedings of the Third Work-</i>		
531	<i>shop on Statistical Machine Translation</i> .		
532	Anthony Chen, Gabriel Stanovsky, Sameer Singh, and		
533	Matt Gardner. 2019. <a href="#">Evaluating question answer-</a>		
534	<a href="#">ing evaluation</a> . In <i>Proceedings of the 2nd Workshop</i>		
535	<i>on Machine Reading for Question Answering</i> , pages		
536	119–124, Hong Kong, China. Association for Com-		
537	putational Linguistics.		
538	Yanran Chen and Steffen Eger. 2023. <a href="#">MENLI: Robust</a>		
539	<a href="#">evaluation metrics from natural language inference</a> .		
540	<i>Transactions of the Association for Computational</i>		
541	<i>Linguistics</i> , 11:804–825.		
542	J. C. Duchi, P. W. Glynn, and H. Namkoong. 2018.		
543	Statistics of robust optimization: A generalized em-		
544	pirical likelihood approach. <i>CoRR</i> , abs/1610.03425.		
545	Alexander R Fabbri, Wojciech Kryściński, Bryan Mc-		
546	Cann, Caiming Xiong, Richard Socher, and Dragomir		
547	Radev. 2021. Summeval: Re-evaluating summariza-		
548	tion evaluation. <i>Transactions of the Association for</i>		
549	<i>Computational Linguistics</i> , 9:391–409.		
550	Marina Fomicheva and Lucia Specia. 2019. <a href="#">Taking MT</a>		
551	<a href="#">evaluation metrics to extremes: Beyond correlation</a>		
552	<a href="#">with human judgments</a> . <i>Computational Linguistics</i> ,		
553	45(3):515–558.		
554	Markus Freitag, Nitika Mathur, Chi kiu Lo, Elefther-		
555	ios Avramidis, Ricardo Rei, Brian Thompson, Tom		
556	Kocmi, Frédéric Blain, Dan Deutsch, Craig Stew-		
557	art, Chrysoula Zerva, Sheila Castilho, Alon Lavie,		
558	and George Foster. 2023. <a href="#">Results of wmt23 met-</a>		
559	<a href="#">rics shared task: Metrics might be guilty but refer-</a>		
560	<a href="#">ences are not innocent</a> . In <i>Proceedings of the Eighth</i>		
561	<i>Conference on Machine Translation</i> , pages 576–626,		
562	Singapore.		
563	Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo,		
564	Craig Stewart, Eleftherios Avramidis, Tom Kocmi,		
565	George Foster, Alon Lavie, and André F. T. Martins.		
566	2022. <a href="#">Results of WMT22 metrics shared task: Stop</a>		
567	<a href="#">using BLEU – neural metrics are better and more</a>		
568	<a href="#">robust</a> . In <i>Proceedings of the Seventh Conference</i>		
569	<i>on Machine Translation (WMT)</i> , pages 46–68, Abu		
570	Dhabi, United Arab Emirates (Hybrid). Association		
571	for Computational Linguistics.		
	F. Maxwell Harper and Joseph A. Konstan. 2015. <a href="#">The</a>	572	
	<a href="#">movielens datasets: History and context</a> . <i>ACM Trans.</i>	573	
	<i>Interact. Intell. Syst.</i> , 5(4).	574	
	Tianxing He, Jingyu Zhang, Tianle Wang, Sachin	575	
	Kumar, Kyunghyun Cho, James Glass, and Yulia	576	
	Tsvetkov. 2023. <a href="#">On the blind spots of model-based</a>	577	
	<a href="#">evaluation metrics for text generation</a> . In <i>Proceed-</i>	578	
	<i>ings of the 61st Annual Meeting of the Association for</i>	579	
	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	580	
	pages 12067–12097, Toronto, Canada. Association	581	
	for Computational Linguistics.	582	
	Yannis Kalantidis, Mert Bulent Sariyildiz, Noe Pion,	583	
	Philippe Weinzaepfel, and Diane Larlus. 2020. Hard	584	
	negative mixing for contrastive learning. In <i>Ad-</i>	585	
	<i>vances in Neural Information Processing Systems</i> ,	586	
	volume 33, pages 21798–21809. Curran Associates,	587	
	Inc.	588	
	Maurice G Kendall. 1938. A new measure of rank	589	
	correlation. <i>Biometrika</i> , 30(1/2):81–93.	590	
	Tom Kocmi, Christian Federmann, Roman Grund-	591	
	kiewicz, Marcin Junczys-Dowmunt, Hitokazu Mat-	592	
	sushita, and Arul Menezes. 2021. <a href="#">To ship or not to</a>	593	
	<a href="#">ship: An extensive evaluation of automatic metrics</a>	594	
	<a href="#">for machine translation</a> . In <i>Proceedings of the Sixth</i>	595	
	<i>Conference on Machine Translation</i> , pages 478–494,	596	
	Online. Association for Computational Linguistics.	597	
	Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Nose-	598	
	worthy, Laurent Charlin, and Joelle Pineau. 2016.	599	
	<a href="#">How NOT to evaluate your dialogue system: An</a>	600	
	<a href="#">empirical study of unsupervised evaluation metrics</a>	601	
	<a href="#">for dialogue response generation</a> . In <i>Proceedings of</i>	602	
	<i>the 2016 Conference on Empirical Methods in Natu-</i>	603	
	<i>ral Language Processing</i> , pages 2122–2132, Austin,	604	
	Texas. Association for Computational Linguistics.	605	
	Nitika Mathur, Johnny Wei, Markus Freitag, Qingsong	606	
	Ma, and Ondřej Bojar. 2020. <a href="#">Results of the WMT20</a>	607	
	<a href="#">metrics shared task</a> . In <i>Proceedings of the Fifth Con-</i>	608	
	<i>ference on Machine Translation</i> , pages 688–725, On-	609	
	line. Association for Computational Linguistics.	610	
	Jekaterina Novikova, Ondřej Dušek, Amanda Cer-	611	
	cas Curry, and Verena Rieser. 2017. <a href="#">Why we need</a>	612	
	<a href="#">new evaluation metrics for NLG</a> . In <i>Proceedings of</i>	613	
	<i>the 2017 Conference on Empirical Methods in Natu-</i>	614	
	<i>ral Language Processing</i> , pages 2241–2252, Copen-	615	
	hagen, Denmark. Association for Computational Lin-	616	
	guistics.	617	
	Vassil Panayotov, Guoguo Chen, Daniel Povey, and San-	618	
	jeev Khudanpur. 2015. <a href="#">Librispeech: An asr corpus</a>	619	
	<a href="#">based on public domain audio books</a> . In <i>2015 IEEE</i>	620	
	<i>International Conference on Acoustics, Speech and</i>	621	
	<i>Signal Processing (ICASSP)</i> , pages 5206–5210.	622	
	Karl Pearson. 1900. X. on the criterion that a given	623	
	system of deviations from the probable in the case	624	
	of a correlated system of variables is such that it	625	
	can be reasonably supposed to have arisen from	626	
	random sampling. <i>The London, Edinburgh, and</i>	627	
	<i>Dublin Philosophical Magazine and Journal of Sci-</i>	628	
	<i>ence</i> , 50(302):157–175.	629	



630	Karl Pearson. 1904. <i>On the theory of contingency and its relation to association and normal correlation</i> . Drapers' Company research memoirs. Cambridge University Press.	In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4344–4355, Florence, Italy. Association for Computational Linguistics.	686 687 688 689
634	Mark Przybocki, Kay Peterson, Sébastien Bronsart, and Gregory Sanders. 2009. The nist 2008 metrics for machine translation challenge—overview, methodology, metrics, and results. <i>Machine Translation</i> , 23:71–103.	Ziang Xiao, Susu Zhang, Vivian Lai, and Q. Vera Liao. 2023. <i>Evaluating evaluation metrics: A framework for analyzing NLG evaluation metrics using measurement theory</i> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 10967–10982, Singapore. Association for Computational Linguistics.	690 691 692 693 694 695 696
639	Ehud Reiter and Anja Belz. 2009. <i>An investigation into the validity of some metrics for automatically evaluating natural language generation systems</i> . <i>Computational Linguistics</i> , 35(4):529–558.	Lifan Yuan, Yangyi Chen, Ganqu Cui, Hongcheng Gao, Fangyuan Zou, Xingyi Cheng, Heng Ji, Zhiyuan Liu, and Maosong Sun. 2024. Revisiting out-of-distribution robustness in nlp: Benchmarks, analysis, and llms evaluations. <i>Advances in Neural Information Processing Systems</i> , 36.	697 698 699 700 701 702
643	Ananya B. Sai, Tanay Dixit, Dev Yashpal Sheth, Sreyas Mohan, and Mitesh M. Khapra. 2021. <i>Perturbation CheckLists for evaluating NLG evaluation metrics</i> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7219–7234, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. In <i>International Conference on Learning Representations</i> .	703 704 705 706
650	Grace S Shieh. 1998. A weighted kendall's tau statistic. <i>Statistics &amp; probability letters</i> , 39(1):17–24.	Kaitlyn Zhou, Su Lin Blodgett, Adam Trischler, Hal Daumé III, Kaheer Suleman, and Alexandra Olteanu. 2022. <i>Deconstructing NLG evaluation: Evaluation practices, assumptions, and their implications</i> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 314–324, Seattle, United States. Association for Computational Linguistics.	707 708 709 710 711 712 713 714 715
652	Miloš Stanojević, Amir Kamran, Philipp Koehn, and Ondřej Bojar. 2015. <i>Results of the WMT15 metrics shared task</i> . In <i>Proceedings of the Tenth Workshop on Statistical Machine Translation</i> , pages 256–273, Lisbon, Portugal. Association for Computational Linguistics.		
658	Amanda Stent, Matthew Marge, and Mohit Singhai. 2005. Evaluating evaluation methods for generation in the presence of variation. In <i>International conference on intelligent text processing and computational linguistics</i> , pages 341–351. Springer.		
663	Tianxiang Sun, Junliang He, Xipeng Qiu, and Xuanjing Huang. 2022. <i>BERTScore is unfair: On social bias in language model-based metrics for text generation</i> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3726–3739, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.		
670	Daniel Valcarce, Alejandro Bellogín, Javier Parapar, and Pablo Castells. 2018. On the robustness and discriminative power of information retrieval metrics for top-n recommendation. In <i>Proceedings of the 12th ACM conference on recommender systems</i> , pages 260–268.		
676	Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson Enrique Yalta Soplín, Jahn Heymann, Matthew Wiesner, Nanxin Chen, Adithya Renduchintala, and Tsubasa Ochiai. 2018. <i>ESPnet: End-to-end speech processing toolkit</i> . In <i>Proceedings of Interspeech</i> , pages 2207–2211.		
683	John Wieting, Taylor Berg-Kirkpatrick, Kevin Gimpel, and Graham Neubig. 2019. <i>Beyond BLEU: Training neural machine translation with semantic similarity</i> .		