



# MENLI: Robust Evaluation Metrics from Natural Language Inference

Yanran Chen<sup>1,2</sup> and Steffen Eger<sup>2</sup>

<sup>1</sup>Technische Universität Darmstadt, Germany

<sup>2</sup>Natural Language Learning Group (NLLG), <https://nl2g.github.io/>

Faculty of Technology, Universität Bielefeld, Germany

yanran.chen@stud.tu-darmstadt.de, steffen.eger@uni-bielefeld.de

## Abstract

Recently proposed BERT-based evaluation metrics for text generation perform well on standard benchmarks but are vulnerable to adversarial attacks, e.g., relating to information correctness. We argue that this stems (in part) from the fact that they are models of semantic similarity. In contrast, we develop evaluation metrics based on *Natural Language Inference* (NLI), which we deem a more appropriate modeling. We design a preference-based adversarial attack framework and show that our NLI based metrics are much more robust to the attacks than the recent BERT-based metrics. On standard benchmarks, our NLI based metrics outperform existing summarization metrics, but perform below SOTA MT metrics. However, when combining existing metrics with our NLI metrics, we obtain both higher adversarial robustness (15%–30%) and higher quality metrics as measured on standard benchmarks (+5% to 30%).

## 1 Introduction

Proper evaluation is key to fields such as machine learning and Natural Language Processing (NLP). Evaluation is particularly challenging for natural language generation (NLG) tasks, as there may be an infinitude of correct solutions (e.g., translations or summaries) for a given source text. While human evaluation is often considered the gold standard, it is slow and costly, thus researchers resort to automatic evaluation. Previously, this was done using simple lexical overlap metrics such as BLEU and ROUGE, but these exhibit low correlations with human judgments, particularly for state-of-the-art NLG systems (Mathur et al., 2020a; Peyrard, 2019). Thus, a popular recent trend is to design automatic evaluation metrics based on large language models such as BERT and its many extensions (Zhang et al., 2020; Zhao et al., 2019; Sellam et al., 2020; Wan et al., 2022).

Nonetheless, these novel metrics also have key limitations. For example, Sai et al. (2021) and Kaster et al. (2021) show that they are not robust to various adversarial attacks including lexical overlap and factuality errors. Taking the currently most popular metric—BERTScore<sup>1</sup>—as an example, this adversarial vulnerability is unsurprising. BERTScore computes the semantic similarity between a reference and a system output (the *candidate*), using a simplified token matching procedure. However, a good candidate is typically not appropriately identified by semantic similarity. For example, a candidate “5 Ukrainian soldiers wounded in Russia” is not an adequate translation of a source corresponding to the reference “50000 Russian soldiers killed in Ukraine”, although the two texts are of course semantically very similar.<sup>2</sup> While there have been many attempts to improve BERTScore using better token matching, e.g., using Word Mover Distance (Zhao et al., 2019; Chen et al., 2020; Colombo et al., 2021), we argue that this line of research is a dead-end, as the underlying model of semantic similarity, originally proposed to address issues of lexical variation in BLEU/ROUGE, is simply not (fully) appropriate.

An intuitively more suitable idea to model evaluation metrics is via *natural language inference* (NLI) (Dagan et al., 2013). For example, in reference-based settings, in which candidates are compared to human references, a candidate is intuitively good if it is *equivalent* to a human reference via the concept of bi-implication. NLI systems are also promising alternatives because

<sup>1</sup>Published in 2020, BERTScore has more than 1700 citations as of March 2023.

<sup>2</sup>That semantic similarity metrics are inherently incapable of identifying this puts them at great risk of being attacked by malicious agents, with serious real-world consequences, as the metrics cannot distinguish between truthful translations and semantically similar but factually incorrect translations.

NLI is one of the most researched upstream tasks in NLP, where a lot of emphasis has been placed on concepts such as biases, generalization and adversarial conditions (Poliak et al., 2018; Utama et al., 2020).

In this paper, we ask whether we can directly use pre-trained NLI models as evaluation metrics, thereby establishing a new paradigm (but with predecessors, as indicated in §2). Our contributions:

- We design: a novel preference-based adversarial test suite for MT and summarization metrics. Our adversarial benchmark does not need human annotators, is suitable for *reference-free* (where the candidate is directly compared to the source text, without human reference) and *reference-based* evaluation, and is challenging: e.g., BLEU, ROUGE, MoverScore, and BERTScore perform below or at random level.
- We explore: (i) how NLI metrics can be induced from existing NLI models; (ii) how they perform on benchmark and adversarial datasets, across (iii) two NLG problems, MT and summarization.
- We show: (iv) NLI metrics perform particularly well in summarization, but below standard metrics in MT. (v) They substantially outperform existing metrics on our adversarial attacks (e.g.,  $\sim 30\%$ – $50\%$  margin over the best unsupervised standard metric in MT). (vi) Combining existing metrics with our NLI metrics yields both better ( $+5\%$ – $30\%$ ) and more robust metrics ( $+15\%$ – $30\%$ ).

We point out that some current metrics already leverage NLI systems—thus, we do not include new information with respect to them—but indirectly and thus (we argue) inadequately: e.g., MoverScore (Zhao et al., 2019) leverages BERT representations fine-tuned on NLI. Mathur et al. (2019) train (pre-BERT) NLI-inspired architectures on MT datasets. In contrast, we show that by *directly* leveraging NLI systems, much better adversarial and standard benchmark performances can be obtained. We call our novel metrics MENLI (**ME**trics from **NLI**).<sup>3</sup>

<sup>3</sup>Code+data: <http://github.com/cyr19/MENLI>.

Concept	Examples
Semantic Similarity	BERTScore, MoverScore, BaryScore, ...
Text Generation	BARTScore, PRISM (Thompson and Post, 2020)
Question Answering	QAEval (Deutsch et al., 2021)
NLI	MENLI

Table 1: Different paradigms for metric induction proposed in recent years.

## 2 Related Work

Our work connects to evaluation metrics and NLI.

**Evaluation Metrics for NLG** In the last few years, researchers have come up with a plethora of different BERT-based metrics for varying tasks and setups: e.g., for MT and summarization, reference-based trained (Sellam et al., 2020; Rei et al., 2020a) and untrained approaches (Zhao et al., 2019; Zhang et al., 2020) have been suggested and the same is true for reference-free setups, where both supervised (Ranasinghe et al., 2020) and unsupervised metrics have been explored (Zhao et al., 2020; Song et al., 2021; Belouadi and Eger, 2023). In our work, we consider both reference-based as well as reference-free metrics. Both setups have important differences: Reference-free setups are more challenging, as they require to compare text in different languages (in MT) or of vastly different lengths (in summarization). On the other hand, they are more ‘resource-efficient’, take humans out-of-the-loop, and promise web-scale evaluation. Both approaches are also different in terms of NLI. For example, while reference-based approaches require equivalence between reference and hypothesis, the concept of equivalence is not always appropriate in reference-free situations (e.g., in summarization, source and summary are intuitively not equivalent; rather, source should entail summary).

To realize metrics, different high-level approaches have been suggested as we outline in Table 1 (e.g., metrics from semantic similarity, from text generation or from question answering). There are also some predecessor works on metrics from NLI which we discuss below.

**Robustness of Evaluation Metrics** has been a central issue of recent interest: Sai et al. (2021) test metrics across several CheckList (Ribeiro et al., 2020) inspired templates, finding that most

common standard metrics are not robust even to simple perturbations. Kaster et al. (2021) probe metrics in an adversarial setting with lexical overlap, finding that they can be fooled by text that has high lexical overlap but low semantic similarity (indicating that the proposed BERT-based metrics are not even good models of semantic similarity). We combine the approaches of Sai et al. (2021) and Kaster et al. (2021): While Sai et al. (2021) use human crowd-workers to evaluate robustness, Kaster et al. (2021) use a simpler preference-based setup, which does not need human annotators. We will also use the preference-based setup, but our attacks are largely inspired by Sai et al. (2021).

More recently (contemporaneously with us and after the first Arxiv submission of our work), several other papers have explored the robustness of recent evaluation metrics. For example, He et al. (2022) develop stress test suites according to potential errors arising from certain choices of metric design and pretrained language models used, showing that metrics are biased towards their underlying models—e.g., BERTScore assigns higher scores to texts generated by the models of the metric itself.<sup>4</sup> Karpinska et al. (2022) explore the sensitivity of MT metrics to errors of different categories (regarding semantics, syntax, and morphology) and severity, using a preference-based setting; they show that recent metrics like BERTScore dramatically outperform lexical overlap-based metrics such as BLEU and ROUGE, mostly obtaining over 95% accuracy in their experiments. Our setups and that of Karpinska et al. (2022) and He et al. (2022) are differentiated by the tasks considered, the preference specifications, the results, and the solutions proposed. Karpinska et al. (2022) only evaluate metrics for MT while we consider both MT and summarization. They design their preferences in such a way that it would seem that recent metrics are quite robust while our more elaborate preferences expose their weak spots much bet-

---

<sup>4</sup>Robustness is also related to model *biases*. For example, Sun et al. (2022) show that BERTScore encodes social biases such as gender biases. And Deutsch et al. (2022) claim that reference-free metrics are inherently biased, which implies that they have unreasonable preferences. Our results show that many current reference-based metrics also have unreasonable preferences. Robustness checks are also related to *explainability* (Leiter et al., 2022; Golovneva et al., 2023) of evaluation metrics as they help to understand metric limitations.

ter. Finally, we propose solutions (e.g., metrics from NLI) to addressing lack of robustness. Like us, He et al. (2022) also consider summarization and MT. Instead of designing preferences, however, they manually introspect how metric scores change as various perturbations are introduced. In this way, they expose blind spots of metrics. As remedies, they suggest to combine heterogeneous metrics to shield against varying blind spots (without performing concrete experiments)—we show that combining metrics with NLI based metrics yields additional robustness.

Finally, Rony et al. (2022) develop RoMe as a robust metric in the context of semantic similarity, fluency and grammatical variability. They evaluate it on an adversarial dataset with five phenomena (entity, adjective and random word replacement; as well as text transformation and passive forms) by correlating against human judgments. Their model is a rather complicated trained metric leveraging semantic and grammatical features—we compare to it in §6.

**NLI** NLI is one of the core upstream tasks in the NLP community. Due to its popularity, NLI has been investigated in-depth, where researchers found that trained models often overfit to low-level statistical cues instead of learning generalizable concepts of logical relationships between sentences (Poliak et al., 2018; Gururangan et al., 2018). As a consequence, many approaches to improve generalization have been investigated (e.g., Belinkov et al., 2019; Utama et al., 2020; Zhou and Bansal, 2020). We argue that a high-quality NLI model would be an excellent candidate for an evaluation metric and explore this in this work.

Like us, Mathur et al. (2019) note the similarity of (MT) evaluation and logical equivalence via NLI. They design supervised MT metrics leveraging different pre-BERT inspired architectures, including one from the NLI community called ESIM (Chen et al., 2017) (which performs on par to an LSTM with attention in their experiments). Thus, in contrast to us, they do not leverage NLI models out-of-the-box as evaluation metrics but only fine-tune an NLI-inspired architecture on human scores from MT. MoverScore (Zhao et al., 2019) fine-tunes BERT on NLI, which leads to better metrics. Thus, they, too, use NLI only indirectly. Dušek and Kasner (2020) use NLI to evaluate hallucinations and omissions in reference-free data-to-text generation scenarios.

	Number error	Negation error
<i>src</i>	Der bilaterale Handel wurde auf über <b>100 Milliarden Dollar</b> im Jahr gesteigert.	Die Wirtschaft der Entwicklungs- und Schwellenländer <b>wird schwach bleiben</b> .
<i>ref</i>	Bilateral trade has increased to more than <b>\$100 billion</b> a year.	Emerging economies <b>will remain weak</b> .
<i>r</i> (google translation of <i>src</i> )	Bilateral trade has increased to over <b>\$100 billion</b> a year.	The economies of developing and emerging countries <b>will remain weak</b> .
<i>cand<sub>para</sub></i>	Bilateral trade has increased to more than <b>one hundred billion dollars</b> a year.	Emerging markets <b>will remain weak</b> .
<i>cand<sub>adv</sub></i> (ref-based)	Bilateral trade has increased to more than <b>\$814 billion</b> a year.	Emerging economies <b>won't remain weak</b> .
<i>cand<sub>adv</sub></i> (ref-free)	Bilateral trade has increased to over <b>\$478 billion</b> a year.	The economies of developing and emerging countries <b>won't remain weak</b> .

Table 2: Examples of our adversarial test suite taken from WMT20<sub>de</sub>. Red words indicate specific adversarial perturbations of the words in green. *cand<sub>adv</sub>*(ref-based) builds on *ref*, whereas *cand<sub>adv</sub>*(ref-free) builds on *r* (indicated by corresponding coloring in the first column). The preferences we query for are given in Eq. (1).

They do not compare to any other metrics and do not consider NLI as a general paradigm for evaluation metrics. While the summarization community uses NLI models for *consistency evaluation* (Fabbri et al., 2021; Laban et al., 2022), to our knowledge, we are the first to verify the usefulness of NLI systems as *general evaluation metrics* against a range of strong competitors, both in standard evaluation and adversarial attack settings.

### 3 Adversarial Setup

Following Sai et al. (2021) and others, we consider an array of adversarial attacks on evaluation metrics—we will give a motivation of our attacks from the perspective of errors committed by real text generation systems below. In contrast to Sai et al. (2021) and similar to the later published work of Karpinska et al. (2022), we implement a preference-based setup, which does not need human annotators. The advantages of the preference-based setup are: (i) lower cost (e.g., no annotation costs), (ii) which can be especially relevant for non-English languages (e.g., in ref-free situations for MT), and (iii) which allows adversarial evaluation at larger scale, yielding more robust estimates of performance. The challenge of the preference setup is to cleverly determine text pairs to compare.

In our design, we use an anchor text (either the reference *ref* or the source *src*), a paraphrase *cand<sub>para</sub>* of the anchor text, and an adversarial text *cand<sub>adv</sub>* which is maximally similar to the anchor text, but contains an adversarial attack. We

expect a good metric *m* to prefer *cand<sub>para</sub>* over *cand<sub>adv</sub>*:

$$\begin{aligned} \text{ref-based} &: m(\text{ref}, \text{cand}_{\text{para}}) > m(\text{ref}, \text{cand}_{\text{adv}}) \\ \text{ref-free} &: m(\text{src}, \text{ref}) > m(\text{src}, \text{cand}_{\text{adv}}) \end{aligned} \quad (1)$$

The outcome of preferences in Eq. (1) depend on how we choose *cand<sub>adv</sub>* and *cand<sub>para</sub>*, which we will describe below. In general, a challenging test suite has *cand<sub>adv</sub>* maximally similar to *ref/src*, but with a key error. In contrast, *cand<sub>para</sub>* should be maximally dissimilar to *ref/src* (e.g., on surface level) but meaning-equivalent. Table 2 illustrates the general structure of our adversarial test suite.

*cand<sub>adv</sub>* To obtain *cand<sub>adv</sub>*, we consider the following attacks (nine regarding information adequacy/correctness in candidates and three regarding text fluency), which we deem (to a large degree) representative for errors in different NLG tasks:

- *Addition*: We randomly add a noun after an existing one and connect them with “and”. For example, “I love dogs” → “I love dogs and cats.”
- *Omission*: We use the framework of Sai et al. (2021) to randomly drop ~1%–20% words in the sentence.
- *Mismatch*: We consider mismatching nouns, verbs, and adjectives, which can lead to misunderstanding of an entity, an action, and the speakers’ emotion, respectively. Following

Chen et al. (2021), we replace a specific word having the POS tag of noun/verb/adjective with another word having the same POS tag randomly selected from our collected words for that POS tag.

- *Negation*: We use the perturbation tool of Ribeiro et al. (2020) to add/remove negations to/from the verb for generating  $cand_{adv}$  with contrary claims.
- *Number error*: We replace all numbers (except for those related to dates) in the sentence with random numbers in the same format (e.g., integer to integer, decimal to decimal).
- *Pronoun error*: We replace all pronouns in the sentence with other ones without causing syntax errors (e.g., ‘‘he’’ to ‘‘she’’ and ‘‘us’’ to ‘‘them’’).
- *Name error*: We use the tool of Ribeiro et al. (2020) to replace exactly one name with a random one of the same gender.
- *Fluency*: We also include three phenomena from Sai et al. (2021) to examine metrics’ robustness against attacks on text fluency: (i) *Jumbling word order*: Randomly shuffle the word order in a sentence. *Spelling error*: Add a typo to a word in a sentence. *Subject-verb disagreement*: Make the subject and verb disagree (e.g., ‘‘He like dogs.’’).

For ref-based metrics, we apply the perturbation templates to  $ref$  to construct  $cand_{adv}$ . In contrast, for ref-free MT metrics, we first translate the source  $src$  using Google Translate to a translation  $r$  and then perturb  $r$  to obtain  $cand_{adv}$ . We introduce  $r$  to increase the similarity of  $cand_{adv}$  to  $src$ ; e.g., we assume that Google Translate translates more literally, i.e., closer to word-by-word translations, than human translators. This may be important to construct challenging test cases, cf. §6 and our above discussion. For ref-free summarization, we apply the perturbation templates to a document  $r$  which is maximally similar to  $src$ ; details follow.

$cand_{para}$  We use different ways to obtain  $cand_{para}$ , because different kinds of paraphrases may yield more/less difficult test cases for metrics. We will analyze this in §6.

In particular, we use data from (1) PAWS (Zhang et al., 2019), (2) PAWS-X (Yang et al., 2019), (3) WMT20-news-commentary-v15 German-to-English (Mathur et al., 2020b) to gen-

dataset	task	ref-based	ref-free	$cand_{para}$	#examples
PAWS <sub>ori</sub>	MT	yes	no	ORI	2,000
PAWS <sub>back</sub>	MT	yes	no	BACK	2,000
XPAWS <sub>x</sub>	MT	yes	yes	ORI	455–474
WMT20 <sub>de</sub>	MT	yes	yes	BACK	200
SE <sub>adv</sub>	SUM	yes	yes	BACK	199

Table 3: Adversarial datasets. ‘‘Yes/no’’ indicates whether the dataset supports ref-based/free adversarial evaluation. ‘‘ORI/BACK’’ denotes whether  $cand_{para}$  (except for *number error*) is from the original datasets or backtranslation. ‘‘#examples’’ refers to the avg. number of examples per phenomenon. XPAWS<sub>x</sub> denotes XPAWS<sub>de/fr/zh/ja</sub>.

erate  $cand_{para}$  for MT evaluation metrics, and (4) SummEval for summarization metrics. A summary with attributes is shown in Table 3.

(1) PAWS contains sentence pairs created by word swapping and backtranslation, labeled as (non-)paraphrases by human raters. From sentence pairs labeled as *paraphrase*, we derive two datasets for ref-based evaluation metrics:

- **PAWS<sub>ori</sub>**: We take the first sentence of a PAWS sentence pair as  $ref$  and the second as  $cand_{para}$ .
- **PAWS<sub>back</sub>**: We use the first sentence of a PAWS sentence pair as  $ref$  and generate  $cand_{para}$  based on  $ref$  using backtranslation (we use German as the pivot language) except for number error, for which we replace the numbers in  $ref$  with the corresponding words, using the Python library `num2words`.

(2) PAWS-X is the multilingual version of PAWS, which includes PAWS sentence pairs in six languages, translated from English PAWS, allowing us to generate test suites for both ref-free and ref-based metrics. We use the first sentence in PAWS-X (e.g., German) as  $src$  and the second sentence with the same ID in English PAWS as  $ref$ . We select the data for two closer language pairs: German-to-English and French-to-English, and two more distant language pairs: Chinese-to-English and Japanese-to-English. Accordingly, we create 4 datasets: **XPAWS<sub>de</sub>**, **XPAWS<sub>fr</sub>**, **XPAWS<sub>zh</sub>**, and **XPAWS<sub>ja</sub>**, each of which contains  $src$  (first sentence of X-PAWS pair in source language),  $ref$  (first sentence of English PAWS

Error	Source	MT hypothesis
Mismatch/verb	关注苏宁易购服务号	<b>Pay attention to (Follow)</b> Suning.com service account
Mismatch/adj.	还不错，玩游戏的画质是真的香	Not bad, the picture quality of playing games is really <b>fragrant (good)</b>
Pronoun/Addition	买给儿子的，他说很好。	Bought it for <b>his (my)</b> son, he said it was good.
Name	当天，美国运输部长赵小兰、美联邦众议员孟昭文以及国际领袖基金会创会会长董继玲等分别在会上发言。	On the same day, US Secretary of Transportation <b>Zhao Xiaolan (Elaine Lan Chao)</b> , US Congressman <b>Meng Zhaowen (Grace Meng)</b> and <b>Dong Jiling (Chiling Tong)</b> , founding president of the International Leaders Foundation, spoke at the meeting respectively.
Omission	I'll <b>review</b> your account, one moment, please.	Ich werde Ihr Konto [...] <b>(überprüfen)</b> , einen Moment bitte.
Mismatch/noun	Listen, I don't want to make <b>my people</b> mad," she said.	„Hör zu, ich will <b>mein Volk (meine Leute)</b> nicht verrückt machen“, sagte sie.
Pronoun	Williams wasn't the only one who received a fine at this year's Wimbledon, though <b>hers</b> was the most costly.	Williams war nicht die einzige, die beim diesjährigen Wimbledon eine Geldstrafe erhielt, obwohl <b>sie (ihre)</b> die teuerste war.

Table 4: Examples of errors in WMT MQM annotations for Chinese-to-English and English-to-German. Red texts are the annotated errors (“[...]” indicates the missing translation) and the green texts in the bracket refer to a more correct translation accordingly; the green texts in source sentences denote the parts being mistranslated or omitted.

pair), and  $cand_{para}$  (second sentence of English PAWS pair).

(3) WMT20-news-commentary-v15 contains sentence pairs of source and human reference. From this, we create **WMT20<sub>de</sub>**, directly taking the source and reference sentences as *src* and *ref*. We obtain  $cand_{para}$  as in the case of **PAWS<sub>back</sub>**.

(4) SummEval (Fabbri et al., 2021) contains documents and references from CNN Daily-Mail (CNNDM) (Hermann et al., 2015), with 10 additional human references. We rank the 11 references using ROUGE-L (Lin, 2004) and use the reference *r* with highest ROUGE score to generate  $cand_{adv}$  for ref-free setting, while the remaining 10 references serve as *ref*. We refer to the adversarial dataset induced from SummEval as **SE<sub>adv</sub>** in the remainder. We obtain  $cand_{para}$  as in the case of **PAWS<sub>back</sub>**.<sup>5</sup>

**Real-world Motivation of Attacks** Modern text generation systems are prone to many of the errors we investigate in this work. For example, Freitag et al. (2021a,b, 2022) show, based on fine-grained human error annotations

<sup>5</sup>As we generate our adversarial test instances fully automatically from backtranslation or automatic tools, they may contain some errors (including upper-/lower-case). For example, we note that in  $cand_{para}$ , “. . . billion dollars” is sometimes incorrectly formulated as “. . . dollars billion”; however, such cases occur only in  $\sim 1\%$  of all test cases for number error, which we argue is still on an acceptable noise level.

(Lommel et al., 2014), that translations generated by state-of-the-art MT models still contain many accuracy-related errors (e.g., addition and omission of information, inappropriately informal pronouns) and sometimes even fluency-related errors (e.g., wrong spelling). Negation handling is also frequently discussed as an issue of modern MT systems (Bentivogli et al., 2016; Sennrich, 2017; Hossain et al., 2020; Tang et al., 2021). In summarization, system summaries are often factually inconsistent with source documents in terms of numbers, named entities and assigning quotations to a particular person, etc. (Falke et al., 2019; Kryscinski et al., 2020; Chen et al., 2021). More generally, *hallucination* (of which addition/mismatches/etc. may be considered special cases) is a particular worrisome limitation of recent large language models (Ji et al., 2022). In Table 4, we show selected system translations from real MT systems with specific errors (following WMT MQM annotations) that are very similar to the ones we consider.<sup>6</sup> The frequency of errors may differ for various source-target language pairs (e.g., depending on their language distance) and formal/informal context. For example, when translating Chinese to English for news, the names are often directly translated to their Pinyin format (see the 4th row) instead of the

<sup>6</sup><https://github.com/google/wmt-mqm-human-evaluation>.

Task		Metrics
MT	ref-based	MoverScore (Zhao et al., 2019), BERTScore (Zhang et al., 2020), BARTScore (Yuan et al., 2021), SentSim (Song et al., 2021), COMET (Rei et al., 2020b), BLEURT (Sellam et al., 2020)
	ref-free	COMET, SentSim, XMoverScore (Zhao et al., 2020)
Summarization	ref-based	BARTScore, DiscoScore (Zhao et al., 2023), MoverScore, BERTScore
	ref-free	BARTScore, SUPERT (Gao et al., 2020)

Table 5: Evaluation metrics explored in this work.

Task		Datasets
MT	segment-level	WMT15-17, WMT20-21
	system-level	WMT20-21
	adversary	<i>ref-based</i> : PAWS <sub>ori/back</sub> , WMT20 <sub>de</sub> , XPAWS <sub>de</sub> ; <i>ref-free</i> : XPAWS <sub>de/fr/zh/ja</sub> , WMT20 <sub>de</sub>
Summarization	summary-level	RealSum (Bhandari et al., 2020)
	system-level	RealSum, SummEval
	adversary	SE <sub>adv</sub> , Rank19 (Falke et al., 2019) (ref-free only)

Table 6: We use the to-English language pairs in WMT15-17 datasets (Stanojević et al., 2015; Bojar et al., 2016, 2017). In segment-level evaluation on WMT20-21 (Mathur et al., 2020b; Freitag et al., 2021a,b), we use the data with MQM scores for zh-en, while in system-level evaluation, we correlate the metrics with DA scores for all to-English language pairs. The datasets for system-level evaluation before WMT20 are skipped, as all metrics mostly get very high correlations on them.

official translations; in contrast, this rarely happens in English-to-German translations. But even for such closely related languages, NLG systems may omit information, or choose wrong pronouns or mismatching nouns, particularly when a word has multiple senses.

## 4 Experimental Setup

### 4.1 Evaluation Metrics

We explore a large array of recent state-of-the-art transformer based metrics, summarized in Table 5. The variants used are briefly introduced below; further details (e.g., model checkpoints and implementation) can be found on our Github.

We report BERTScore F1 employing a RoBERTa-large model. For MoverScore, we use the unigram variant with a BERT-base model fine-tuned on MNLI (Williams et al., 2018). We use two variants of BARTScore (Precision and F1) for ref-based MT and summarization and BARTScore-FN (FN stands for Faithfulness) for ref-free summarization. We consider two variants of XMoverScore with different remapping strategies for multilingual embeddings (CLP, UMD) and two variants of SentSim with different word matching paradigms (BERTScore, WMD). We report the DiscoScore variant with feature ‘Focus Frequency’.

### 4.2 Datasets & Evaluation Protocol

We summarize our used datasets in Table 6. To evaluate the metrics’ robustness under **adversarial conditions**, we use the datasets introduced in §3 and additionally Rank19 (Falke et al., 2019) (only for ref-free summarization), which contains examples composed of documents paired with one correct and one incorrect candidate summary with real-world factuality errors. In general, we check the metrics’ preference between the two candidates and calculate *accuracy*: the relative frequency that the metrics correctly choose among the two alternatives.

On **MT** standard benchmarks, we evaluate the metrics on both *segment-level* (where we correlate metrics scores to human judgments for individual sentences/segments in the datasets) and *system-level* (where we correlate the average metric scores to the average human scores over the segments generated by each system), using Pearson correlation as the performance indicator. On SummEval for **summarization**, we compute Kendall correlation with system-level human judgements on four criteria: *coherence*, *consistency*, *fluency* and *relevance* (we apply two aggregation methods for the multi-reference setting, *max* and *mean*). We calculate Pearson correlation with both summary-level (analogous to

segment-level in MT) and system-level *LitePyramids* (Shapira et al., 2019) human ratings in RealSumm.

### 4.3 NLI as a Metric

NLI systems yield probability distributions over *Entailment*, *Contradiction*, and *Neutral*. We denote the probability values as  $e$ ,  $c$ , and  $n$ , where  $e + c + n = 1$  and  $e, c, n \geq 0$ . We first determine how to leverage the three values as NLI metrics.

To do so, we evaluate **five simple formulas** of their arithmetic combination in a heuristic way: (1)  $e$ , (2)  $-c$ , (3)  $e-n$ , (4)  $e-c$ , and (5)  $e-n-2c$ , and inspect their effect in **three directions**, which correspond to the entailment directions implication, reverse implication and bi-implication: (i)  $ref/src \rightarrow cand$ , where  $ref$  or  $src$  act as premise and  $cand$  as hypothesis; (ii)  $ref/src \leftarrow cand$ , where  $cand$  acts as premise and  $ref$  or  $src$  act as hypothesis; and (iii)  $ref/src \leftrightarrow cand$ , as arithmetic average over the two above cases.

For example, to obtain  $e-n$  from  $ref/src \leftrightarrow cand$ , we first average the three probability scores over direction  $ref/src \rightarrow cand$  and  $ref/src \leftarrow cand$ , then calculate  $e-n$  based on the averaged scores. We only consider direction  $src \rightarrow cand$  for ref-free summarization, since hypothesis does not need to entail source document. The various selections of the formulas and directions result in 15 pooling strategies for NLI-based metrics.

**NLI Systems** We explore both monolingual and cross-lingual NLI-based metrics. For each setup, we choose two NLI models, which are obtained from Hugging Face or fine-tuning by ourselves.

For **monolingual NLI metrics**, we choose (1) a RoBERTa-large model (Liu et al., 2019) fine-tuned on SNLI (Bowman et al., 2015), MNLI, Fever (Nie et al., 2019) and ANLI (Nie et al., 2020) by Nie et al. (2020) and (2) a DeBERTa-large model fine-tuned by He et al. (2021), using MNLI. We denote the NLI metrics induced from these two models as  $NLI-R$  and  $NLI-D$ . They will be used for ref-based MT evaluation, and both ref-based and -free summarization evaluation tasks. Note that, while  $NLI-R$  has been fine-tuned on adversarial NLI (ANLI), which has been shown to increase robustness on (for example) negation and numerical reasoning,  $NLI-D$  has not been trained on ANLI. **Cross-lingual NLI metrics** should handle premises and hypotheses in different languages, so we select the multilingual versions of the under-

(a) Reference-based					
	e	-c	e-n	e-c	e-n-2c
$ref \rightarrow cand$	3+0	3+0		2+0	
$ref \leftarrow cand$					
$ref \leftrightarrow cand$	0+4		0+3	0+1	0+2
(b) Reference-free					
	e	-c	e-n	e-c	e-n-2c
$src \rightarrow cand$		2+0			
$src \leftarrow cand$	0+1		0+2		
$src \leftrightarrow cand$	0+1		4+6	4+0	

Table 7: Winning frequency of different pooling strategies for NLI metrics on adversarial (first entry) and MT datasets (second entry). We only show non-zero entries.

lying models of  $NLI-R/NLI-D$ . (1) We fine-tune a XLM-RoBERTa-base model (Conneau et al., 2019), using the datasets for fine-tuning  $NLI-R$  as well as XNLI dataset (Conneau et al., 2018). (2) We select an mDeBERTa-base model fine-tuned on MNLI and XNLI. We denote the corresponding cross-lingual NLI metrics as  $XNLI-R$  and  $XNLI-D$ .

## 5 Experiment Results

Before outlining our main results in §5.1 (MT) and §5.2 (summarization), we first discuss good pooling strategies for NLI metrics.

**Pooling Strategy** We determine the pooling strategy for NLI metrics in **MT evaluation** from (1) the accuracy on the adversarial datasets and (2) the correlation with human judgements on the standard (segment-level) MT datasets. We leverage the *winning frequency* of the pooling strategies to choose the best one; a strategy wins if it works best for an NLI metric among all 15 strategies. Overall, we find that the simple formula  $e$  from the direction  $src/ref \leftrightarrow cand$  is a good choice which works well for both standard and adversarial benchmarks, even though slightly better formulas could be chosen in selected subsetings (e.g., ref-based vs. ref-free evaluation), see Table 7 for examples.

For **summarization**, the situation is slightly more complex: (1)  $e-c$  from direction  $ref \leftarrow cand$  performs best for ref-based NLI metrics; (2)  $-c$  from direction  $src \rightarrow cand$  is the best strategy for



	Adv.				MT			
	ref-based		ref-free		ref-based		ref-free	
	all	adeq.	all	adeq.	seg	sys	seg	sys
Supervised								
COMET	67.4	67.0	76.8	74.5	0.676	0.808	0.620	0.698
BLEURT	74.8	79.8	–	–	0.708	0.807	–	–
Unsupervised								
sentBLEU	32.9	27.2	–	–	0.380	0.757	–	–
Rouge	34.3	28.7	–	–	0.425	0.774	–	–
MoverScore	48.3	46.9	–	–	0.567	<b>0.806</b>	–	–
XMoverS(UMD)			74.5	71.7	–	–	0.400	0.672
XMoverS(CLP)	–	–	73.8	70.9	–	–	0.422	<b>0.673</b>
BERTS	65.3	60.9	–	–	<b>0.620</b>	0.799	–	–
BARTS-P	67.4	64.2	–	–	0.587	0.761	–	–
BARTS-F	78.4	77.8	–	–	0.593	0.802	–	–
SentS(BERTS)	68.1	67.8	62.7	65.5	0.612	0.401	0.421	−0.021
SentS(WMD)	62.1	61.9	63.0	65.8	0.607	–	<b>0.427</b>	–
NLI-based								
X(NLI)-R	85.0	92.1	70.5	75.8	0.451	0.756	0.221	0.335
X(NLI)-D	<b>86.6</b>	<b>92.3</b>	<b>79.3</b>	<b>85.8</b>	0.439	0.770	0.149	0.581

Table 8: Pearson correlation with human judgments in WMT and accuracy (%) on our adversarial datasets, averaged over datasets. The performance of ref-based COMET is averaged over WMT20<sub>de</sub> and XPAWS<sub>de</sub>, since it also requires source texts as input. In bold: best results among all unsupervised metrics including the NLI-based metrics.

ref-free NLI metrics. Thus, we compare NLI metrics adopting these strategies with classic metrics.

Even though we only looked at global aggregate statistics, we still observe that our method of identifying the pooling strategies above leveraged the data on which we will later evaluate the NLI metrics. To avoid leaking information from the test set, we evaluate NLI metrics on each dataset with the pooling strategy selected from the remaining datasets for that task in §6.

## 5.1 Machine Translation

### 5.1.1 Adversarial Evaluation

We now compare our NLI metrics with the best pooling strategy to our baseline metrics under adversarial conditions.

From Table 8 (columns ‘‘Adv.’’), we observe that in the **ref-based** setup: (1) NLI metrics outperform the great majority of metrics by a huge margin: over 85% vs. 32%–78% (all phenomena) and 92% vs. 27%–80% (adequacy phenomena only) on average. (2) Further, the two NLI metrics perform similarly. In the **ref-free** setup, the best cross-lingual NLI metric (XNLI-D) is still most robust under our attacks. However, NLI metrics do not as substantially outperform the other met-

rics as in the ref-based setup. A potential reason is that the cross-lingual NLI models underperform compared to the monolingual setup (the preferences we query for in the reference-free setup may also play a role). Nevertheless, when excluding the fluency-related phenomena from the adversarial datasets, XNLI-D is still on average 10 points better than the best standard metric, COMET (86% vs. 75%).

Moreover, our results reveal that: (1) most standard metrics are particularly incapable of detecting *name error*, *number error*, and *pronoun error* (~29%–70%); (2) standard metrics, especially BLEURT and COMET, are most competitive regarding *omission*, *addition*, and *jumbling* (~80%–100%); (3) NLI metrics are sub-optimal for fluency attacks (mostly at random level), especially the reference-free NLI metrics; and (4) NLI metrics are much better at *name error*, *negation*, *number error*, *pronoun error*, and *adj. mismatch* than most of the other metrics, especially ref-based (>90% vs. ~10%–80%), as shown in Figure 1.

Our observations are inconsistent with Karpinska et al. (2022), where the state-of-the-art MT metrics mostly obtain >95% accuracy in the preference-based evaluation. The reason is that our test suites are much more difficult for the evaluation metrics because we challenge them by lexical overlap between source/reference and candidate sentences during attacks: Metrics must choose between high lexical overlap adversarial candidates (with key errors) over low lexical overlap paraphrases. In contrast, in Karpinska et al. (2022), metrics are challenged to assign correct preferences for  $\text{score}(\text{ref}, t)$  vs.  $\text{score}(\text{ref}, t')$  where  $t$  is a candidate and  $t'$  the perturbed candidate. This is a much easier comparison because neither are  $\text{ref}$  and  $t$  maximally dissimilar (but meaning equivalent) nor are  $\text{ref}$  and  $t'$  maximally similar. This is an important lesson: *How to design the adversarial preferences may critically affect the assessment of whether recent metrics are robust or not.*

### 5.1.2 Standard Benchmarks

**Ref-based** We give average results over all datasets in Table 8 (columns ‘MT’; individual results are available in our Github). For **segment-level** evaluation, we observe: (1) trained metrics (COMET and BLEURT) substantially outperform the others, with average performance

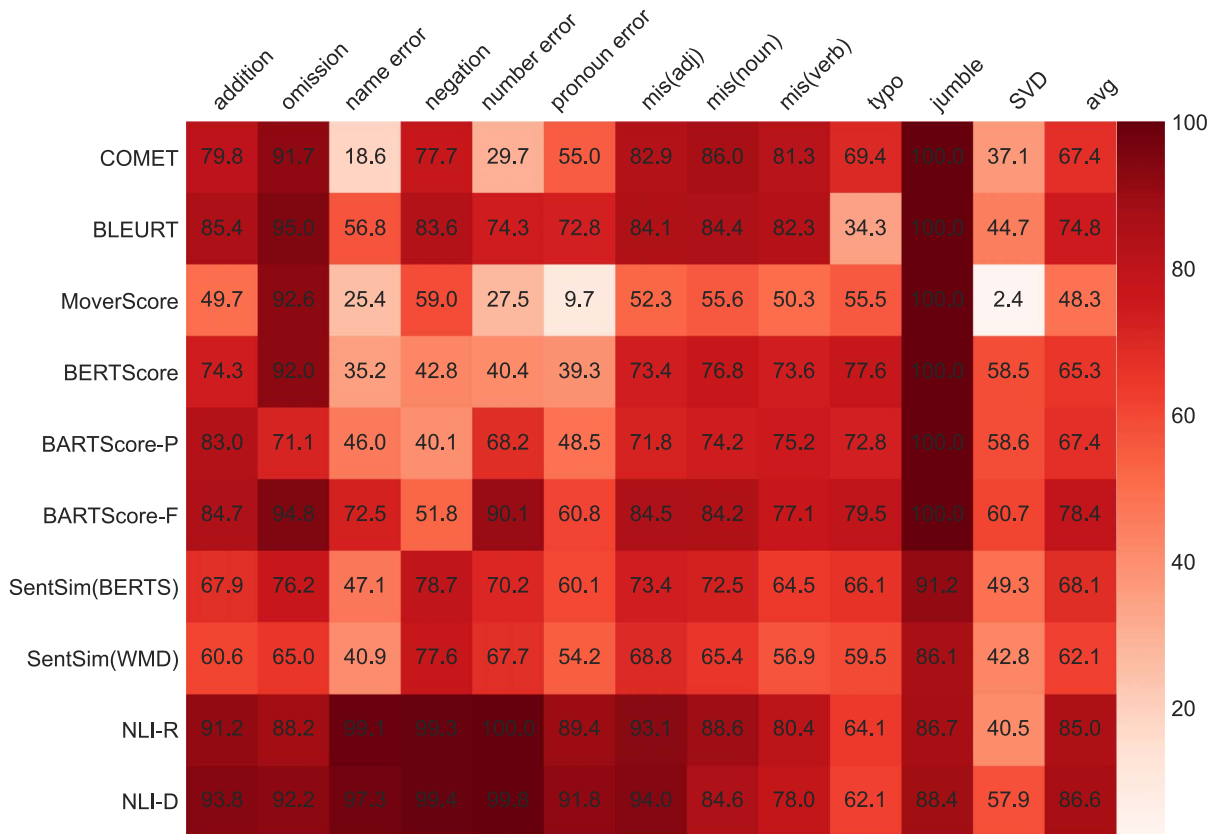


Figure 1: Average accuracy (values in each block) of all metrics per phenomenon over the adversarial datasets for ref-based MT evaluation. Darker color indicates higher accuracy and vice versa.

of  $\sim 0.7$  Pearson. (2) Unsupervised SOTA metrics have average correlation of  $\sim 0.6$  Pearson, BERTScore is the best among them. (3) Our NLI-based metrics are not competitive, with correlations of  $\sim 0.45$  Pearson. When correlating with **system-level** human judgments, NLI metrics still underperform most of the SOTA metrics, but the margin is much smaller.

**Ref-free** Trained metrics also dominate in **segment-level** evaluation ( $>0.6$  Pearson), whereas the two NLI-based metrics perform much worse than the others (0.15-0.22 Pearson). Nevertheless,  $\text{XNLI-D}$  performs on par with COMET and better than the others on WMT20 at **system-level**.

Overall, we conclude that our NLI metrics are not competitive with state-of-the-art evaluation metrics on standard MT datasets, especially at segment-level and ref-free.

### 5.1.3 Combined Metrics

Observing that NLI metrics are strong on adversarial setups, but comparatively weaker in standard

evaluation, we examine *how to get more robust metrics which also perform well on standard benchmarks*. To do so, we take the weighted average of NLI and classical metrics:

$$C = w_{\text{nli}} \cdot N + (1 - w_{\text{nli}}) \cdot M \quad (2)$$

where  $w_{\text{nli}} \in [0, 1]$  is the weight for NLI metric  $N$  and  $M$  is a classical metric. Before combination, we rescale  $M$  and  $N$  to  $[0, 1]$ , using min-max normalization.

We illustrate the performance of the combined evaluation metrics with  $(\text{X})\text{NLI-R}$  on both adversarial and standard benchmarks (segment-level) in Figure 2; the results for  $(\text{X})\text{NLI-D}$  and for system-level are similar. The  $x$ -axis denotes the average accuracy over the adversarial datasets, while  $y$ -axis is the average Pearson correlation over the standard benchmarks (MT datasets). Each dot in each graph shows the value  $C(w_{\text{nli}})$  for a specific weight  $w_{\text{nli}}$ . As seen from Figure 2, the graphs show an intriguing concave curvature. In standard MT evaluation, the combination boosts the metric

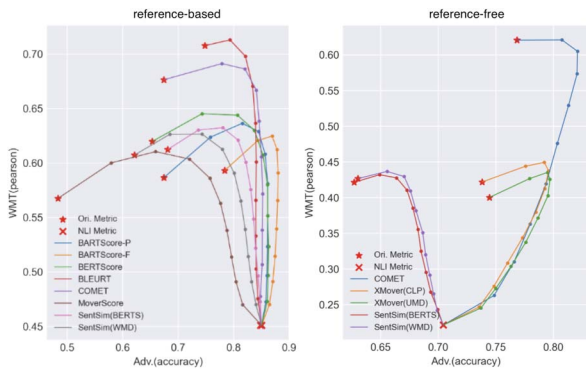


Figure 2: Accuracy on adversarial datasets and Pearson correlation with **segment-level** human judgements in WMT datasets of combined metrics with (X)NLI-R, averaged over datasets. The points on each path from the original metric to the NLI metric indicate  $w_{\text{nli}} = 0, 0.1, \dots, 1$ . The purple line denoting the combination with ref-based COMET ends at another point since the corresponding adversarial performance is averaged over the 2 adversarial datasets containing source texts.

performance when  $w_{\text{nli}}$  is small (from 0.1 to 0.4) in virtually all cases. We then see a *simultaneous* increase of adversarial robustness and quality on standard benchmarks. In **ref-based** setup, e.g., for  $w_{\text{nli}} = 0.2$ , we observe: (1) MoverScore and BARTScore-P improve most, with  $\sim 8\%$  (from 0.57/0.59 to 0.61/0.64 Pearson, respectively) and 21%–36% improvements on adversarial datasets (from 48%/67% to 66%/82% accuracy on average). (2) The best unsupervised metric on segment-level MT, BERTScore, increases  $\sim 4\%$  Pearson on standard benchmarks and  $\sim 24\%$  accuracy on adversarial datasets. (3) The most robust untrained metric, BARTScore-F, improves about  $\sim 11\%$  in robustness, whereas its performance on standard benchmarks also rises  $\sim 5\%$ . (4) The improvements on MT for trained metrics are smaller compared to those untrained metrics, with COMET improving only 1.5% and BLEURT even becoming worse with the choice  $w_{\text{nli}} = 0.2$ . However, their performance in defending adversarial attacks still improves  $\sim 10\%$ – $20\%$ . In **ref-free** setups, all metrics improve  $\sim 6\%$ – $7\%$  on adversarial datasets. Such setting only substantially boosts XMoverScore’s performance on standard benchmarks, with  $\sim 6\%$ – $9\%$ .

We summarize the improvements for all combinations in Figure 3(a), which are averages over all experiments considered here. We can observe that the line denoting improvements on standard

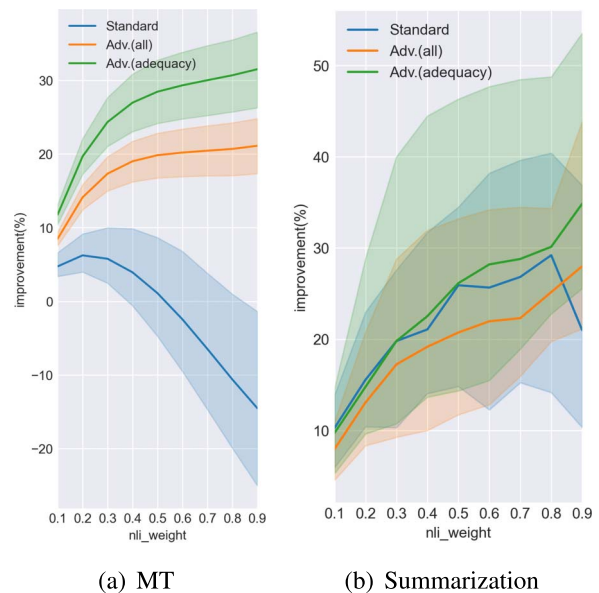


Figure 3: Improvements of all metrics on standard benchmarks and adversarial datasets for  $w_{\text{nli}} = 0.1, \dots, 0.9$ , averaged over all experiments. We show 95% confidence interval.

benchmarks peaks at  $w_{\text{nli}} = 0.2$ , and the average improvements are positive when  $w_{\text{nli}} \leq 0.5$ . Further, on the adversarial datasets, the improvement monotonously increases with  $w_{\text{nli}}$  and the gain is a concave function of  $w_{\text{nli}}$  which saturates as  $w_{\text{nli}}$  becomes larger. The sweet spots are  $w_{\text{nli}} \in [0.2, 0.3]$ , which leads to 5%–6% improvement on standard benchmarks and 14%–16% improvement in adversarial robustness on average. When excluding the fluency phenomena from the adversarial datasets, the combined metrics consistently gain larger improvements in adversarial robustness, with 20%–24% improvements at the sweet spots.

## 5.2 Summarization

**Evaluation** As Table 9 shows, similar to MT evaluation, NLI-based metrics exhibit much stronger robustness under adversarial conditions (our best NLI metrics have at least  $\sim 8$  points higher accuracy than the best standard metrics; right-most columns). The difference is that the vanilla NLI metrics are now also comparably effective to the SOTA metrics on standard benchmarks. For instance, in **ref-based** setup, NLI-D with *max* aggregation beats all metrics except for DiscoScore with *mean* on SummEval and both NLI metrics highly correlate with system-level human ratings in RealSumm (above

(a) Reference-based														
metric	SummEval										RealSumm		Adv.	
	coherence		consistency		fluency		relevance		avg		litePyr		SE <sub>adv</sub>	
	mean	max	mean	max	mean	max	mean	max	mean	max	sum	sys	all	adeq.
BLEU	0.294	0.279	0.044	-0.029	0.244	0.229	0.397	0.382	0.245	0.215	0.480	0.124	0.182	0.109
Rouge	0.191	0.176	0.088	-0.279	-0.037	-0.081	0.118	0.103	0.090	-0.020	0.540	0.457	0.185	0.117
MoverS	0.206	0.324	0.456	0.103	0.421	0.362	0.368	0.515	0.363	0.326	<b>0.585</b>	0.501	0.287	0.251
BERTS	0.618	0.618	0.221	0.044	0.273	0.185	0.603	0.515	0.429	0.340	0.574	0.380	0.598	0.574
BARTS-P	0.485	0.441	0.176	-0.044	0.376	0.185	0.500	0.368	0.385	0.237	0.478	0.531	0.697	0.692
BARTS-F	0.515	<b>0.647</b>	0.206	0.250	0.317	0.450	0.529	<b>0.632</b>	0.392	0.495	0.583	0.687	0.788	0.792
DiscoS	<b>0.676</b>	0.279	0.279	0.676	0.539	0.554	<b>0.632</b>	0.353	<b>0.532</b>	0.466	-0.199	-0.066	0.334	0.294
NLI-based														
NLI-R	0.147	0.074	0.632	0.676	0.494	0.450	0.279	0.206	0.388	0.352	0.525	<b>0.856</b>	<b>0.864</b>	<b>0.905</b>
NLI-D	0.250	0.265	<b>0.706</b>	<b>0.750</b>	<b>0.568</b>	<b>0.613</b>	0.471	0.397	0.499	<b>0.506</b>	0.489	0.840	0.806	0.843

(b) Reference-free														
metric	SummEval					RealSumm		Adv.						
	coherence		consistency	fluency	relevance	avg	litePyr		SE <sub>adv</sub>		Rank19			
	summary	system	all	adeq.	avg	summary	system	all	adeq.	avg				
BARTS-FN	<b>0.735</b>		0.132	0.391	<b>0.662</b>	<b>0.480</b>	0.178	-0.023	0.427	0.389	0.796	0.612		
SUPERT	0.147		0.603	<b>0.465</b>	0.279	0.374	<b>0.522</b>	0.626	0.296	0.273	0.668	0.482		
NLI-based														
NLI-R	0.221		0.235	0.391	0.500	0.337	0.300	<b>0.688</b>	<b>0.720</b>	<b>0.722</b>	0.866	<b>0.793</b>		
NLI-D	0.162		<b>0.647</b>	0.332	0.324	0.366	-0.076	0.568	0.624	0.629	<b>0.885</b>	0.755		

Table 9: Kendall correlation with system-level human judgments in SummEval. Pearson correlation with summary/system-level litePyramid in RealSumm. Accuracy on adversarial benchmarks, averaged over phenomena in SE<sub>adv</sub>. We bold the best performance on each criterion. ‘‘max/mean’’ denotes the aggregation method used for multi-reference setting in ref-based evaluation on SummEval.

0.8 Pearson), where most standard metrics obtain only 0.5–0.7 Pearson correlations. When considering all evaluation dimensions of SummEval and RealSumm, NLI-D outperforms all other metrics, followed by NLI-R. Besides, we observe that NLI metrics correlate much better with human judgments regarding *consistency* and (somewhat surprisingly) *fluency* in SummEval compared to the other metrics. For the **ref-free** setup, BARTScore-FN performs best on SummEval—it outperforms the other metrics by above 0.1 Kendall on average. However, it does not correlate well with both summary-level and system-level human judgments in RealSumm. NLI metrics are comparable or better than standard metrics on system-level. For example, NLI-R performs best among the examined metrics and is about 0.06 Pearson better than the best standard metric (SUPERT) on system-level in RealSumm. Nevertheless, reference-free NLI metrics also perform worse than the reference-based ones as in MT; an explicit bottleneck for the two NLI metrics is that they were only trained on NLI data with short sentences, but reference-free

summarization evaluation requires metrics to deal with source documents which contain many more sentences.

**Combined Metrics** In Figure 3(b), we summarize the median improvements of combined summarization metrics (the median smooths some outliers). In contrast to MT, the combination brings almost equal benefits to performance of standard metrics on standard and adversarial benchmarks concerning only adequacy—we again observe a decrease in improvements on adversarial datasets when adding our fluency phenomena. We identify a best  $w_{nli}$ , namely, 0.8, with which the standard metrics gain about 25%–30% improvements in *both* types of performances (adversarial and standard).

## 6 Discussion & Analysis

**Selected Failure Cases of Metrics:** Table 10 shows selected failure cases of four popular metrics (BERTScore, BARTScore, BLEURT, COMET), where the NLI metrics are correct in

<i>ref</i>	<i>cand<sub>para</sub></i>	<i>cand<sub>adv</sub></i>	<i>score<sub>para</sub></i> : <i>score<sub>adv</sub></i> (standard metric)	<i>score<sub>para</sub></i> : <i>score<sub>adv</sub></i> (NLI-R)	<b>error</b>
BERTScore					
Although President George W. Bush says <b>he</b> believes in markets, in this case <b>he</b> has called for voluntary action.	Although President George W. Bush says <b>he</b> believes in markets, <b>he</b> has demanded voluntary action in this case.	Although President George W. Bush says <b>she</b> believes in markets, in this case <b>she</b> has called for voluntary action.	0.980: 0.982	0.951: 0.000	<i>Pronoun</i>
BARTScore-F					
<b>Reagan</b> and I were nonetheless able to create a reservoir of constructive spirit through constant outreach and face-to-face interaction.	Nevertheless, <b>Reagan</b> and I were able to create a constructive climate through constant contact and personal interaction.	<b>Nicole</b> and I were nonetheless able to create a reservoir of constructive spirit through constant outreach and face-to-face interaction.	-2.104: -1.527	0.943: 0.002	<i>Name</i>
BLEURT					
In 2012, when Freedom House downgraded <b>Mali</b> to “not free,” engagement declined by 7%.	In 2012, when Freedom House classified <b>Mali</b> as unfree, the engagement fell by 7 percent.	In 2012, when Freedom House downgraded <b>Melissa</b> to “not free,” engagement declined by 7%.	0.787: 0.834	0.983: 0.030	<i>Name</i>
This leads to heavy deforestation and lethal indoor air pollution, which kills <b>1.3 million</b> people each year.	This leads to heavy Deforestation and lethal indoor air pollution, which kills <b>one point three million</b> people each year.	This leads to heavy Deforestation and lethal indoor air pollution, which kills <b>6.9 million</b> people each year.	0.682: 0.767	0.783: 0.000	<i>Num</i>
COMET					
Who serves as president of the United States <b>is</b> critically important for Mexicans.	Anyone who serves as President of the United States <b>is</b> crucial to Mexicans.	Who serves as president of the United States <b>is not</b> critically important for Mexicans.	1.067: 1.086	0.974: 0.044	<i>Negation</i>

Table 10: Sample instances in adversarial datasets where standard metrics failed while NLI-R succeeded; ref-based setup. In the 4th and 5th columns, we show [*score assigned to cand<sub>para</sub>*]: [*score assigned to cand<sub>adv</sub>*] by standard metrics and NLI-R, respectively; robust metrics should give *cand<sub>para</sub>* higher scores. Green bold texts indicate the anchor words/phrases to be perturbed and the red ones in *cand<sub>adv</sub>* refer to the corresponding perturbed texts.

each case. In the examples, BERTScore prefers text with the wrong gendered pronoun over a legitimate paraphrase and even trained metrics like BLEURT fail on severe name changes such as “Melissa” (a person name) vs. “Mali” (a country name). Leveraging more subtle cases (e.g., mismatches based on wrong word senses instead of random mismatches with the same POS or replacing names with names of the same ‘type’) would likely constitute even harder test cases for future metrics.

**No Metric is Good Everywhere:** Across distinct dimensions, different metrics perform differently, indicating that they capture varying aspects. For example, NLI metrics are not so good on fluency adversarial attacks, e.g., typos. This may be unsurprising, given that fluency is a low-level phenomenon while NLI concerns high-level logical relationships between sentences (some fluency phenomena would best be treated by switch-

ing to a lower-level representation space, such as character-level [Vu et al., 2022]; this could seamlessly be integrated in existing NLI models). The NLI metrics are also weaker concerning segment-level MT evaluation on standard benchmarks. However, NLI metrics alone perform surprisingly well: In ref-based MT, they win on 7 out of 19 dimensions (12 adversarial phenomena and 7 standard datasets, evaluated segment- and system-level), only beaten by BLEURT (8 wins); ref-free, they win 5 out of 19 dimensions, second only to COMET (11 wins). In ref-based summarization, they are clearly ahead of all standard metrics, winning not only 8 out of 12 adversarial dimensions, but also system-level LitePyramid, consistency and fluency (thus, 11 out of 18 wins), clearly ahead of BARTScore-P (4 of 18); ref-free, they are also best and win 13 out of 18 dimensions. The best overall metrics, measured as average performance over standard and adversarial datasets, always include NLI: for ref-based MT,

this is BLEURT+0.2×NLI-R, for ref-free MT, it is COMET+0.3×NLI-D. For summarization, NLI-R alone and combined with BARTScore-F perform best on average.

**Rescaling:** The min-max normalization we used (a standard technique for normalizing data in machine learning, typically applied to input features) for metric combination requires batch processing. It is necessary to account for the different ranges of metrics, e.g., some metrics take negative values. An alternative would include to enforce more formal constraints on evaluation metrics, i.e., that they should take outputs in [0,1]. When applying our combined metrics in practice, one could also replace them by surrogate metrics trained on the outputs of the original combined metrics or simply take the min-max values inferred from the datasets already evaluated on—the larger these datasets the more reliably are min and max estimated.

**Sensitivity to  $w_{\text{nli}}$ :** Having different weights  $w_{\text{nli}}$  for different tasks is undesirable, because it requires considering each task individually. However, in our experiments, we found that all small  $w_{\text{nli}}$  (below 0.5) yield good performances and are thus safe choices: They increase adversarial robustness and also lead to better metrics on standard benchmarks.

**Adversarial Performance vs. Standard Performance:** From our experiments, it might seem that adversarial and standard performance are anti-correlated: A metric with higher adversarial performance may have lower performance on standard benchmarks and vice versa. While this would not necessarily be a major surprise as adversarial conditions oftentimes test phenomena that are otherwise not represented in standard benchmarks (Niven and Kao, 2019), a statistical analysis reveals that standard performance generally *positively* correlates to the adversarial performance in our case, consistent with our earlier argument that existing NLG systems in the real world do commit similar errors as we check for. To do so, we first convert the metrics’ standard performance to rankings for each performance category (e.g., ref-based/-free segment/system-level MT performance, performance on SummEval/RealSumm), then we correlate the ranking-based standard performance to the corresponding adversarial performance rankings, obtaining 0.37 Spearman.

When excluding NLI metrics, the correlation increases to 0.60.

**The Choice of  $\text{cand}_{\text{para}}$  Matters:** As indicated in §3, we speculate that a good adversarial setting maximizes (surface) dissimilarity between *ref* and  $\text{cand}_{\text{para}}$  (which can better trick the metrics). To investigate, we compute the normalized edit distance between *ref* and  $\text{cand}_{\text{para}}$ ,<sup>7</sup> a larger edit distance means a greater dissimilarity. If our assumption is true, then larger edit distances represent harder test cases for the metrics. We find: (1) the average edit distance for the test cases where the metrics fail to defend against the adversarial attacks is 0.01–0.6 larger than that for where they succeed, averaged over metrics; (2) for PAWS<sub>back</sub> and PAWS<sub>ori</sub> (both induced from PAWS) where the  $\text{cand}_{\text{para}}$  are obtained in different ways, all metrics achieve 0.02–0.15 lower accuracy on PAWS<sub>ori</sub>, which has 0.46 larger average edit distance than PAWS<sub>back</sub>, in turn. Both findings confirm our above assumption. In addition, we observe that NLI metrics have the smallest difference between the edit distances for failure and success cases (0.01–0.26) as well as that between the accuracy on PAWS<sub>back</sub> and PAWS<sub>ori</sub> (0.02) among all evaluation metrics. This implies that they are least affected by surface overlap and instead better consider the logical relationship between sentences. This is what makes them attractive as evaluation metrics.

**The Choice of  $\text{cand}_{\text{adv}}$  Matters, Too:** We evaluate on one complex attack combining *Number error* with *Negation* which increases the difference between *ref* and  $\text{cand}_{\text{adv}}$  based on the test cases for *Number error* in WMT20<sub>de</sub>. The accuracy increases by an average of 0.28 over all metrics. This confirms our assumption that maximizing the (surface) similarity between *ref* and  $\text{cand}_{\text{adv}}$  (but with key errors) leads to harder test suites and vice versa.

**Ensemble with NLI Metrics Are More Effective:** We compare the ensembles with NLI metrics to ensembles with standard metrics, i.e.,  $w \cdot A + (1 - w) \cdot M$ , where  $A$  is a fixed standard metric and  $M$  is any of the remaining metrics. To do so, we combine standard metrics with the rest metrics for each category of MT/summarization and ref-based/-free setting. We take the arithmetic

<sup>7</sup>Ref-free, the edit distance between  $r$  and *ref* is considered.



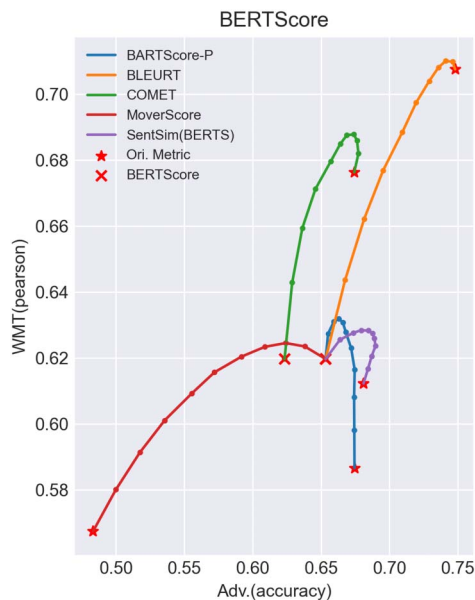


Figure 4: Accuracy on adversarial datasets and Pearson correlation with segment-level human judgements in WMT datasets of combined metrics with BERTScore, averaged over datasets. The green line denoting the combination with COMET ends at another point since the corresponding adversarial performance is only averaged over the 2 adversarial datasets containing source texts.

average of the accuracy on adversarial benchmarks and correlations on standard benchmarks as the overall metric performance here. We calculate the mean/maximal improvement of ensembles to the original metric  $M$  over  $w \in [0.1, 0.9]$  and observe: (i) While the ensembles with standard metrics are better for ref-free MT metrics because cross-lingual NLI metrics perform very poorly in our experiments, (ii) the monolingual NLI metrics lead to much better ensembles—17/15 points larger mean/max improvement—compared to the standard metrics. (iii) Overall, the ensembles with NLI metrics yield 10/7 points larger mean/max improvement in overall performance than with standard metrics (averaged over all 4 tasks: ref-based/-free MT/summarization). Thus, (monolingual) NLI metrics have unique properties, compared to standard metrics, making them attractive in ensembles.

To illustrate, Figure 4 shows ensembles with BERTScore. These show minor or no improvements on standard benchmarks and also mixed (often negative) results for adversarial robustness.

**SummaCZS and Falsesum:** In §5, we applied NLI systems on whole input texts, not taking into

account the multi-sentence nature of source texts and outputs, especially in summarization.

To remedy the mismatch between the granularities of the training data of NLI models and the input data of summarization evaluation, i.e., sentence- vs. document-level, Laban et al. (2022) propose both supervised and unsupervised NLI-based summarization metrics for inconsistency detection. We test their unsupervised variant (**SummaCZS**),<sup>8</sup> which segments documents into sentence units and aggregates scores between pairs of sentences, with the underlying model of NLI-R. However, SummaCZS does not consistently outperform NLI-R across all datasets; in contrast, NLI-R performs much better in our adversarial test compared to SummaCZS (72% vs. 53%). Besides, to match the training data of NLI models with the task of factual inconsistency detection in summarization, Utama et al. (2022) introduce an augmented NLI dataset with task-oriented examples based on CNNDM—**FalseSum**; we evaluate three Roberta-large models finetuned on it and MNLI. Similar to SummaCZS, this also does not always yield better performance compared to simple NLI metrics ( $\sim 55\%$ – $68\%$  vs. 72% on adversarial datasets). Overall, both approaches work well on SummEval, but not so well on RealSumm and our adversarial benchmark.

**Choice of Pooling Strategy:** To examine the issue of data leakage discussed in §5, we now evaluate the NLI metrics on each dataset with the pooling strategy selected from the remaining datasets (excluding the one for evaluation) based on winning frequency. For example, for the segment-level MT evaluation on WMT15, we choose the pooling strategy which wins most times on all MT datasets (including all standard datasets for both segment/system-level evaluation and the adversarial datasets) except for WMT15. We observe that this change in pooling strategy induction results in minor performance variation:  $-1.9\%$  for segment-level evaluation,  $+0.8\%$  for system-level evaluation, and  $-0.7\%$  for adversarial evaluation. For summarization, as only one direction—i.e.,  $src \rightarrow cand$ —is considered for ref-free NLI metrics, we separately select the pooling strategy for ref-based and

<sup>8</sup>We do not compare to the supervised one as it is trained on a consistency dataset for summarization task, for a fairer comparison.

ref-free NLI metrics. Overall, we have no performance change for the ref-free setting and  $-3.6\%$  performance on average over all five criteria (correlations on SummEval with max/mean aggregation, summary/system-level correlations on RealSumm, and accuracy on  $SE_{adv}$ ) ref-based. Thus, the changes are again minor.

**Comparison to RoMe:** As the authors of RoMe did not publish their adversarial dataset, we compare RoMe’s performance with our metrics on one of our adversarial datasets, WMT20<sub>de</sub>, instead. RoMe has an average accuracy of 43%, with  $> 90\%$  accuracy only on the phenomena *SVD* and *omission*, which are the easiest for most standard metrics. In contrast, our NLI metrics have above 80% average accuracy. As RoMe does not evaluate on MT or summarization, we also evaluate our NLI metrics on one (randomly chosen) data-to-text generation dataset used in Rony et al. (2022)—BAGEL (Mairesse et al., 2010). RoMe and our NLI metrics perform on par here ( $\sim 0.23$  Spearman’s  $\rho$ ). Overall, this seems to imply that simple NLI models taken out of the box are better and more robust metrics than a specially trained approach such as RoMe.

## 7 Concluding Remarks

In this work, we explored NLI as a *general* paradigm for evaluation metrics. We showed that NLI metrics yield adversarial robustness, and are also strong—though not always state-of-the-art—when it comes to standard metric evaluation benchmarks. By linearly interpolating established (BERT-based) metrics with our NLI metrics, we obtained high-quality metrics along both axes: adversarial robustness and standard benchmarks, with substantial gains over recent BERT-based metrics.

A potential reason why NLI based metrics perform subpar on some standard benchmarks (especially in MT) is the training data mismatch, i.e., typical NLI datasets contain many artificial sentences of the type “A girl is playing on a piano”. A further limitation is that cross-lingual NLI models are not yet high-quality enough and that most current NLI models are sentence-level, not document-level—with a few recent exceptions (Yin et al., 2021). Once these limitations of NLI are overcome, we believe that even better performances from NLI based metrics can be expected,

which, we believe, is one of the most promising directions for future high-quality and robust evaluation metric design. Future work should also consider NLI metrics for other text generation tasks; the NLI paradigm looks especially promising for tasks that require comparison with human references, which oftentimes involve the concept of logical equivalence.

## Acknowledgments

We thank Zuojun Shi for conducting initial experiments related to this paper as part of her Bachelor thesis at TU Darmstadt. We appreciate the reviewers and editors from TACL for their time, effort, and greatly helpful comments. We also thankfully acknowledge support from the BMBF via the grant “Metrics4NLG”. Steffen Eger is financed by DFG grant EG 375/5–1.

## References

- Yonatan Belinkov, Adam Poliak, Stuart Shieber, Benjamin Van Durme, and Alexander Rush. 2019. On adversarial removal of hypothesis-only bias in natural language inference. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (\*SEM 2019)*, pages 256–262, Minneapolis, Minnesota. Association for Computational Linguistics. <https://doi.org/10.18653/v1/S19-1028>
- Jonas Belouadi and Steffen Eger. 2023. Uscore: An effective approach to fully unsupervised evaluation metrics for machine translation. In *EACL*.
- Luisa Bentivogli, Arianna Bisazza, Mauro Cettolo, and Marcello Federico. 2016. Neural versus phrase-based machine translation quality: A case study. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 257–267, Austin, Texas. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D16-1025>
- Manik Bhandari, Pranav Narayan Gour, Atabak Ashfaq, Pengfei Liu, and Graham Neubig. 2020. Re-evaluating evaluation in text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language*



- Processing (EMNLP)*, pages 9347–9359, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.emnlp-main.751>
- Ondřej Bojar, Yvette Graham, and Amir Kamran. 2017. Results of the WMT17 metrics shared task. In *Proceedings of the Second Conference on Machine Translation*, pages 489–513, Copenhagen, Denmark. Association for Computational Linguistics. <https://doi.org/10.18653/v1/W17-4755>
- 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 Papers*, pages 199–231, Berlin, Germany. Association for Computational Linguistics. <https://doi.org/10.18653/v1/W16-2302>
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D15-1075>
- Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Si Wei, Hui Jiang, and Diana Inkpen. 2017. Enhanced LSTM for natural language inference. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1657–1668, Vancouver, Canada. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P17-1152>
- Xi Chen, Nan Ding, Tomer Levinboim, and Radu Soricut. 2020. Improving text generation evaluation with batch centering and tempered word mover distance. In *Proceedings of the First Workshop on Evaluation and Comparison of NLP Systems*, pages 51–59, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.eval4nlp-1.6>
- Yiran Chen, Pengfei Liu, and Xipeng Qiu. 2021. Are factuality checkers reliable? Adversarial meta-evaluation of factuality in summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2082–2095, Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.findings-emnlp.179>
- Pierre Colombo, Guillaume Staerman, Chloé Clavel, and Pablo Piantanida. 2021. Automatic text evaluation through the lens of Wasserstein barycenters. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10450–10466, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.emnlp-main.817>
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116. <https://doi.org/10.18653/v1/2020.acl-main.747>
- Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. Xnli: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D18-1269>
- Ido Dagan, Dan Roth, Mark Sammons, and Fabio Massimo Zanzotto. 2013. *Recognizing Textual Entailment: Models and Applications*, Synthesis Lectures on Human Language Technologies, Morgan & Claypool Publishers. <https://doi.org/10.1007/978-3-031-02151-0>
- Daniel Deutsch, Tania Bedrax-Weiss, and Dan Roth. 2021. Towards question-answering as an automatic metric for evaluating the content quality of a summary. *Transactions of the Association for Computational Linguistics*, 9:774–789. [https://doi.org/10.1162/tacl\\_a\\_00397](https://doi.org/10.1162/tacl_a_00397)
- Daniel Deutsch, Rotem Dror, and Dan Roth. 2022. On the limitations of reference-free evaluations

- of generated text. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10960–10977, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ondřej Dušek and Zdeněk Kasner. 2020. Evaluating semantic accuracy of data-to-text generation with natural language inference. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 131–137, Dublin, Ireland. Association for Computational Linguistics.
- Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. Summeval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409. [https://doi.org/10.1162/tacl\\_a\\_00373](https://doi.org/10.1162/tacl_a_00373)
- Tobias Falke, Leonardo F. R. Ribeiro, Prasetya Ajie Utama, Ido Dagan, and Iryna Gurevych. 2019. Ranking generated summaries by correctness: An interesting but challenging application for natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2214–2220. <https://doi.org/10.18653/v1/P19-1213>
- 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. [https://doi.org/10.1162/tacl\\_a\\_00437](https://doi.org/10.1162/tacl_a_00437)
- 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.
- Yang Gao, Wei Zhao, and Steffen Eger. 2020. SUPERT: Towards new frontiers in unsupervised evaluation metrics for multi-document summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1347–1354, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.124>
- Olga Golovneva, Moya Peng Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2023. ROSCOE: A suite of metrics for scoring step-by-step reasoning. In *The Eleventh International Conference on Learning Representations*.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics. <https://doi.org/10.18653/v1/N18-2017>
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. DeBERTa: Decoding-enhanced BERT with disentangled attention. In *International Conference on Learning Representations*.
- Tianxing He, Jingyu Zhang, Tianle Wang, Sachin Kumar, Kyunghyun Cho, James Glass, and Yulia Tsvetkov. 2022. On the blind spots of model-based evaluation metrics for text generation. *arXiv preprint arXiv:2212.10020*.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. *Advances in Neural Information Processing Systems*, 28.
- Md Mosharaf Hossain, Antonios Anastasopoulos, Eduardo Blanco, and Alexis Palmer. 2020.

- It's not a non-issue: Negation as a source of error in machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3869–3885, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.findings-emnlp.345>
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2022. Survey of hallucination in natural language generation. *ACM Computing Surveys*.
- Marzena Karpinska, Nishant Raj, Katherine Thai, Yixiao Song, Ankita Gupta, and Mohit Iyyer. 2022. DEMETR: Diagnosing evaluation metrics for translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9540–9561, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Marvin Kaster, Wei Zhao, and Steffen Eger. 2021. Global explainability of BERT-based evaluation metrics by disentangling along linguistic factors. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8912–8925, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.emnlp-main.701>
- Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the factual consistency of abstractive text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9332–9346, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.emnlp-main.750>
- Philippe Laban, Tobias Schnabel, Paul N. Bennett, and Marti A. Hearst. 2022. SummaC: Revisiting NLI-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177. [https://doi.org/10.1162/tacl\\_a\\_00453](https://doi.org/10.1162/tacl_a_00453)
- Christoph Leiter, Piyawat Lertvittayakumjorn, M. Fomicheva, Wei Zhao, Yang Gao, and Steffen Eger. 2022. Towards explainable evaluation metrics for natural language generation. *ArXiv*, abs/2203.11131.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Arle Lommel, Aljoscha Burchardt, and Hans Uszkoreit. 2014. Multidimensional quality metrics (mqm): A framework for declaring and describing translation quality metrics. *Tradumàtica: Tecnologies de la Traducció*, 0:455–463. <https://doi.org/10.5565/rev/tradumatica.77>
- François Mairesse, Milica Gašić, Filip Jurčiček, Simon Keizer, Blaise Thomson, Kai Yu, and Steve Young. 2010. Phrase-based statistical language generation using graphical models and active learning. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 1552–1561, Uppsala, Sweden. Association for Computational Linguistics.
- Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2019. Putting evaluation in context: Contextual embeddings improve machine translation evaluation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2799–2808, Florence, Italy. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1269>
- Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020a. Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4984–4997, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.448>
- Nitika Mathur, Johnny Wei, Markus Freitag, Qingsong Ma, and Ondřej Bojar. 2020b. Results of the WMT20 metrics shared task. In

- Proceedings of the Fifth Conference on Machine Translation*, pages 688–725, Online. Association for Computational Linguistics.
- Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. Combining fact extraction and verification with neural semantic matching networks. In *Association for the Advancement of Artificial Intelligence (AAAI)*. <https://doi.org/10.1609/aaai.v33i01.33016859>
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.441>
- Timothy Niven and Hung-Yu Kao. 2019. Probing neural network comprehension of natural language arguments. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4658–4664, Florence, Italy. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1459>
- Maxime Peyrard. 2019. Studying summarization evaluation metrics in the appropriate scoring range. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5093–5100, Florence, Italy. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1502>
- Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 180–191, New Orleans, Louisiana. Association for Computational Linguistics. <https://doi.org/10.18653/v1/S18-2023>
- Tharindu Ranasinghe, Constantin Orasan, and Ruslan Mitkov. 2020. TransQuest: Translation quality estimation with cross-lingual transformers. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5070–5081, Barcelona, Spain (Online). International Committee on Computational Linguistics. <https://doi.org/10.18653/v1/2020.coling-main.445>
- Ricardo Rei, Craig Stewart, Ana C. Farinha, and Alon Lavie. 2020a. 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. <https://doi.org/10.18653/v1/2020.emnlp-main.213>
- Ricardo Rei, Craig Stewart, Ana C. Farinha, and Alon Lavie. 2020b. Unbabel’s participation in the WMT20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 911–920, Online. Association for Computational Linguistics.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.442>
- Md Rashad Al Hasan Rony, Liubov Kovriguina, Debanjan Chaudhuri, Ricardo Usbeck, and Jens Lehmann. 2022. RoMe: A robust metric for evaluating natural language generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5645–5657, Dublin, Ireland. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.acl-long.387>
- Ananya B. Sai, Tanay Dixit, Dev Yashpal Sheth, Sreyas Mohan, and Mitesh M. Khapra. 2021. Perturbation checklists for evaluating NLG evaluation metrics. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of ACL*. <https://doi.org/10.18653/v1/2020.acl-main.704>
- Rico Sennrich. 2017. How grammatical is character-level neural machine translation?

- Assessing MT quality with contrastive translation pairs. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 376–382, Valencia, Spain. Association for Computational Linguistics. <https://doi.org/10.18653/v1/E17-2060>
- Ori Shapira, David Gabay, Yang Gao, Hadar Ronen, Ramakanth Pasunuru, Mohit Bansal, Yael Amsterdamer, and Ido Dagan. 2019. Crowdsourcing lightweight pyramids for manual summary evaluation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 682–687, Minneapolis, Minnesota. Association for Computational Linguistics. <https://doi.org/10.18653/v1/N19-1072>
- Yurun Song, Junchen Zhao, and Lucia Specia. 2021. Sentsim: Crosslingual semantic evaluation of machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3143–3156. <https://doi.org/10.18653/v1/2021.naacl-main.252>
- 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. <https://doi.org/10.18653/v1/W15-3031>
- Tianxiang Sun, Junliang He, Xipeng Qiu, and Xuanjing Huang. 2022. BERTScore is unfair: On social bias in language model-based metrics for text generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3726–3739, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Gongbo Tang, Philipp Rönchen, Rico Sennrich, and Joakim Nivre. 2021. Revisiting negation in neural machine translation. *Transactions of the Association for Computational Linguistics*, 9:740–755. [https://doi.org/10.1162/tacl\\_a\\_00395](https://doi.org/10.1162/tacl_a_00395)
- 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. <https://doi.org/10.18653/v1/2020.emnlp-main.8>
- Prasetya Ajie Utama, Joshua Bambrick, Nafise Sadat Moosavi, and Iryna Gurevych. 2022. Falsesum: Generating document-level nli examples for recognizing factual inconsistency in summarization. <https://doi.org/10.48550/arXiv.2205.06009>
- Prasetya Ajie Utama, Nafise Sadat Moosavi, and Iryna Gurevych. 2020. Mind the trade-off: Debiasing NLU models without degrading the in-distribution performance. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8717–8729, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.770>
- Doan Nam Long Vu, Nafise Sadat Moosavi, and Steffen Eger. 2022. Layer or representation space: What makes BERT-based evaluation metrics robust? In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3401–3411, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Yu Wan, Dayiheng Liu, Baosong Yang, Haibo Zhang, Boxing Chen, Derek Wong, and Lidia Chao. 2022. UniTE: Unified translation evaluation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8117–8127, Dublin, Ireland. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2022.acl-long.558>
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics. <https://doi.org/10.18653/v1/N18-1101>

- Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. 2019. PAWS-X: A cross-lingual adversarial dataset for paraphrase identification. In *Proceedings of EMNLP*. <https://doi.org/10.18653/v1/D19-1382>
- Wenpeng Yin, Dragomir Radev, and Caiming Xiong. 2021. DocNLI: A large-scale dataset for document-level natural language inference. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4913–4922, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2021.findings-acl.435>
- Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. BARTscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems*, 34:27263–27277.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTscore: Evaluating text generation with BERT. In *International Conference on Learning Representations*.
- Yuan Zhang, Jason Baldridge, and Luheng He. 2019. PAWS: Paraphrase adversaries from word scrambling. In *Proceedings of NAACL*.
- Wei Zhao, Goran Glavaš, Maxime Peyrard, Yang Gao, Robert West, and Steffen Eger. 2020. On the limitations of cross-lingual encoders as exposed by reference-free machine translation evaluation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1656–1671, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.151>
- Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. MoverScore: Text generation evaluating with contextualized embeddings and earth mover distance. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 563–578, Hong Kong, China. Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-1053>
- Wei Zhao, Michael Strube, and Steffen Eger. 2023. Discoscore: Evaluating text generation with BERT and discourse coherence. In *EACL*.
- Xiang Zhou and Mohit Bansal. 2020. Towards robustifying NLI models against lexical dataset biases. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8759–8771, Online. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.acl-main.773>