

# Mitigating the Impact of Reference Quality on Evaluation of Summarization Systems with Reference-Free Metrics

Anonymous ACL submission

## Abstract

Automatic metrics are used as proxies to evaluate abstractive summarization systems when human annotations are too expensive. To be useful, these metrics should be fine-grained, show a high correlation with human annotations, and ideally be independent of reference quality; however, most standard evaluation metrics for summarization are reference-based, and existing reference-free metrics correlate poorly with relevance, especially on summaries of longer documents. In this paper, we introduce a reference-free metric that correlates well with human evaluated relevance, while being very cheap to compute. We show that this metric can also be used along reference-based metrics to improve their robustness in low quality reference settings.

## 1 Introduction

Given an input source, an abstractive summarization system should output a summary that is short, relevant, readable and consistent with the source. To reflect this, fine-grained human evaluations are split into different scores (Fabbri et al., 2021), such as fluency, faithfulness (sometimes called factual consistency), coherence and relevance. Fluency measures the linguistic quality of individual sentences, *eg* if they contain no grammatical errors. Coherence gauges if sentences in a summary are well-organized and well-structured. Faithfulness, or factual consistency, considers factual alignment between a summary and the source. Relevance is the measure of whether a summary contains the main ideas from the source.

Automatic summarization metrics are intended to capture one or multiple of these qualities (Zhu and Bhat, 2020; Vasilyev and Bohannon, 2021a), and used as a proxy to evaluate summarization systems when human annotations are too expensive.

These metrics can be compared on their different attributes such as the reliance on one or multiple

references, the cost of inference (Wu et al., 2024), the dataset-agnosticism (Faysse et al., 2023) and their correlations with human judgment at system-level (Deutsch et al., 2022) or summary-level.

In this work, we introduce a new reference-free metric that intends to capture the relevance of machine summaries using  $n$ -gram importance weighting. We rate  $n$ -grams of the source documents relative to how much semantic meaning they express, as measured by *tf-idf* (Sparck Jones, 1972), and score summaries according to their weighted lexical overlap with these  $n$ -grams.

We show that this metric is complementary to other metrics and can be mixed with reference-based metrics to alleviate their sensitivity to noisy and low quality references.

## 2 Related Work

### 2.1 Reference-based evaluation

Lexical overlap based metrics such as ROUGE (Lin, 2004), BLEU (Papineni et al., 2002) and chrF (Popović, 2015), or pretrained language model based metrics such as BERTScore (Zhang et al., 2019) and BARTScore (Yuan et al., 2021), are the standard way of evaluating abstractive summarization systems. However, these metrics rely on gold standard reference summaries that can be costly, noisy, or missing altogether. We discuss some of the limits of these methods in section 3.

### 2.2 LLM-as-a-Judge evaluation

Large Language Models (LLM) are able to perform many tasks with good results, even in a few-shot or even zero-shot fashion. Recently, LLMs are even used to evaluate natural language generation tasks in replacement for human evaluation. LLM-as-a-Judge show useful properties as an evaluation metric, for instance Faysse et al. (2023) illustrated using GPT-4 that it can be highly correlated with human judgement, format and task agnostic and

comparable across tasks. Zheng et al. (2023) describe limitations of LLM-as-a-Judge, including position, verbosity and self-enhancement biases as well as poor performance at grading math or reasoning tasks. Other limitations are expressed by Kim et al. (2023) targeting proprietary LLMs such as GPT-4 for their closed source nature, uncontrolled versioning, and their high costs. Prometheus 2 (Kim et al., 2024) is designed for evaluating language models and show high correlations with proprietary LLMs and human evaluations. Besides, its open-source nature mitigates some of the aforementioned issues. Liu et al. (2023) suggest that LLMs aligned from human feedback overfit to reference-less human evaluation of summaries, which they observed to be biased towards longer summaries and to suffer from low inter-annotator agreement.

### 2.3 Reference-free evaluation

Metrics designed to evaluate summaries without reference are useful when no gold reference are available, or when the property they intend to capture does not need a reference to be conveniently estimated.

GRUEN (Zhu and Bhat, 2020) aims at estimating the linguistic quality of a given summary by taking into account the grammaticality, non-redundancy, focus, structure and coherence of a summary. These attributes can be assessed without relying on a reference summary or even the source document. ESTIME (Vasilyev and Bohannon, 2021a) is evaluating the inconsistencies between the summary and the source by counting the mismatched embeddings out of the hidden layer of a pretrained language model.

Liu et al. (2023) observed that reference-free human evaluations have a very low correlation with reference-based human evaluations, and tend to be biased towards different types of systems.

### 2.4 Evaluating Summarization of Long Documents

Trained metrics usually generalize poorly to out-of-distribution tasks (Koh et al., 2022), and often cannot handle long contexts. In the long document summarization setting, Koh et al. (2022) showed that most automatic metrics correlate poorly with human judged relevance and factual consistency scores. Wu et al. (2024) use an extract-then-evaluate method to reduce the size of the long source document used as a reference for evaluation of factual consistency and relevance with LLM-as-

a-Judge. They find that it both lowers the cost of evaluation, and improve the correlation with human judgement.

### 3 Limits of reference-based evaluation

Lexical overlap scores such as BLEU or ROUGE work under the implicit assumption that reference summaries are mostly extractive and contain no errors.

This assumption is challenged by a study conducted by Maynez et al. (2020) on hallucinated content in abstractive summaries. In human written summaries from the XSum dataset, 76.9% of the gold references were found to have at least one hallucinated word.

Summarization methods can trade abstractiveness for faithfulness, creating a faithfulness-abstractiveness tradeoff curve that was illustrated and studied by Ladhak et al. (2022). They show that some metrics are more sensitive to the summary abstractiveness than others.

In the context of translations, *translationese* refers to source language artifacts found in both human and machine translations. This phenomenon is similar to extractive segments in summaries, as it is an artifact of the source document that can be mitigated by paraphrasing. Freitag et al. (2020) demonstrated that reference translations in machine translation datasets tend to exhibit this *translationese* language. They addressed this by creating new references through paraphrasing the existing ones. When tested, systems produced much lower BLEU scores with the paraphrased references compared to the *translationese* ones, but the correlation with human judgment was higher. They observed that with *translationese* references, the  $n$ -grams with the highest match rates resulted from translations adhering to the source sentence structure. In contrast, using the paraphrased references, the most-matched  $n$ -grams were related to the semantic meaning of the sentence.

Following a *translationese* - extractiveness analogy, we assume that with highly extractive references, the most matched  $n$ -grams between proposed and reference summaries are artifacts of the extractiveness of the summaries. More abstractive references will yield much lower ROUGE scores, but might correlate better with human judgement.

We propose to use  $n$ -gram importance weighting methods, such as *tf-idf* (Sparck Jones, 1972) or *bm-25* (Robertson and Jones, 1976), to extract the

$n$ -grams expressing most of the semantic meaning of the source document. We believe that these  $n$ -grams should appear in relevant summaries, and are not artifacts of extractiveness.

#### 4 Proposed Metric

Let  $W_{t,d,D}$  be the importance of a  $n$ -gram  $t$  in a document  $d$  from a corpus  $D$ . The importance score  $W_{t,d,D}$  is defined as

$$W_{t,d,D} = \begin{cases} \tanh\left(\frac{w_{t,d,D}}{r_{t,d,D}}\right), & \text{if } t \in d \\ 0, & \text{otherwise,} \end{cases}$$

$w_{t,d,D}$  is an importance score obtained through word importance scoring methods (such as *tf-idf* and *bm-25*). The associated importance rank of the  $n$ -gram in the document is referred as  $r_{t,d,D}$ .

Given a proposed summary  $\hat{s}$  of a document  $d \in D$ , we compute the metric:

$$m(\hat{s}, d, D) = \frac{\alpha_{\hat{s},d,D}}{N_{d,D}} \sum_{t \in \hat{s}} W_{t,d,D}$$

With  $N_{d,D}$  the upper ceiling of the sum of weights, used to normalize the score:  $N_{d,D} = \sum_{t \in d} W_{t,d,D}$ .

By construction this score will be maximum for a summary consisting of the full document. To alleviate this issue, we penalize longer summaries by multiplying with a function  $f$  accounting for the length of the summary  $|\hat{s}|$  and the length of the document  $|d|$ :  $\alpha_{\hat{s},d} = f(|\hat{s}|, |d|)$ <sup>1</sup>.

We observe that this length penalty not only resolves the issue related to the scoring of entire documents but also shows a stronger correlation with human judgment at the system level.

Another issue is that it is relatively straightforward to devise a trivial heuristic that achieves a perfect score by employing the same  $n$ -gram importance weighting method to generate an extractive summary, with access to the full corpus. We do not consider this point to be a substantial issue, as such heuristic will result in a low score on metrics that measure other aspects of a summary, such as fluency or faithfulness.

#### 5 Experiments

For our experiments, we work with different datasets of human evaluation of summarization systems.

SummEval (Fabbri et al., 2021) contains human evaluations for 23 systems, each with 100 summaries of news article from the CNN/DailyMail dataset. Coherence, consistency, fluency and relevance are evaluated by experts and crowd-source workers. ArXiv and GovReport (Koh et al., 2022) contain annotations for 12 summarization systems, evaluated on 18 long documents for each dataset. Human evaluators rated the factual consistency and the relevance of the machine summaries. RoSE (Liu et al., 2023) is a benchmark consisting of 12 summarization systems evaluated on 100 news article from CNN/DailyMail. Each summary is annotated with different protocols, we are using the reference-based and reference-free human evaluations.

We describe the choice of settings for our metric in Appendix A, which takes into account system-level correlations on the four datasets, as well as the range of values taken by the metric.

##### 5.1 System-level correlation scaling with number of summaries

According to Deutsch et al. (2022), system-level correlations are usually inconsistent with the practical use of automatic evaluation metrics. When computing these correlations to evaluate systems, usually only the subset of summaries judged by humans is used. However automatic metrics can be computed on summaries outside of this subset to give better estimates. With our reference-free metric we could go one step further and evaluate the systems on more documents, without reference summaries. As illustrated by Deutsch et al. (2022), testing with more examples will also narrow down the confidence intervals of the evaluated scores, making it easier to compare systems.

Figure 1 show an increase of the system-level correlation with human evaluated relevance when using more examples for each system.

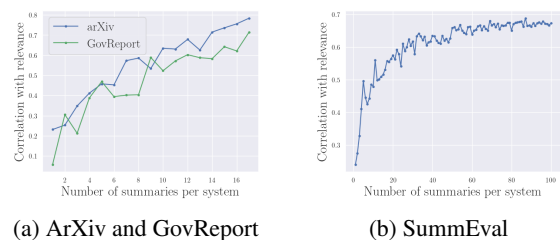


Figure 1: Scaling of system-level correlations with human judgement for our metric, depending on the number of summaries used for evaluation

<sup>1</sup>The choice for  $f$  is illustrated in Appendix A, Figure 4

## 5.2 Robustness to noisy references

Reference-based metrics such as ROUGE-1 are sensitive to the quality of the references. To evaluate the robustness of ROUGE-1 to noisy references, we gradually replace random reference summaries with the first three sentences of the source document and compute the resulting system-level correlations. Results, averaged over 20 random draws, are reported in Figure 2. Our metric is not sensitive to altered references by construction, contrary to ROUGE-1. When mixed with it, it improves the robustness of ROUGE-1 to low quality references. This aspect is beneficial in settings where the quality of the reference summaries is unknown or variable, for instance with web-crawled datasets.

We observe a surprising behaviour of ROUGE-1 where its system-level correlations are improved on average when introduced to alterations in the range of 0 to 30% of the samples of the SummEval dataset. This behaviour does not seem to translate to the mix of ROUGE-1 with our metric.

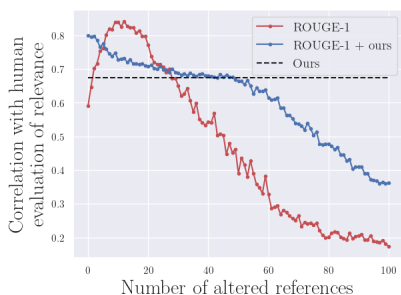


Figure 2: System-level correlation depending on the number of altered references in SummEval

## 5.3 Complementarity with other automatic metrics

We report the pairwise complementarity between each pair of metric<sup>2</sup> on SummEval in Figure 3, following Colombo et al. (2023). We observe that our metric has a high complementarity with most other metrics, especially with ESTIME, intended to capture the factual consistency. Our metric also has a high complementarity with ROUGE and chrF scores, which are also based on lexical overlap, meaning that they capture different features of the evaluated summaries.

In Table 1 we report the system-level Spearman correlations on SummEval using our metric, other

<sup>2</sup>we use the `evaluate` implementation of ROUGE, chrF and BERTScore and official implementations of GRUEN and ESTIME

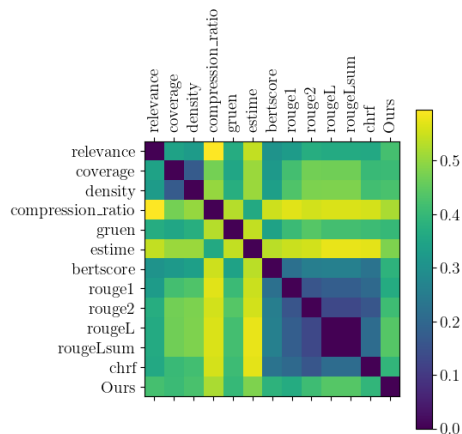


Figure 3: Complementarity between metrics on SummEval

Table 1: System-level correlations of mixes of metrics

Metric	SummEval Spearman
ROUGE-1	0.59
ROUGE-2	0.61
ROUGE-L	0.47
chrF	<b>0.75</b>
BERTScore	0.40
ESTIME	-0.45
GRUEN	0.59
<b>Ours</b>	<b>0.73</b>
<b>Ours + ROUGE-1</b>	<b>0.86</b>
<b>Ours + chrF</b>	0.74
<b>Ours + BERTScore</b>	0.77
<b>Ours - ESTIME</b>	0.76
<b>Ours + GRUEN</b>	<b>0.83</b>

metrics, and mixing our metric with other metrics.

## 6 Conclusion and future works

In this work, we introduce a new reference-free metric based on importance-weighted  $n$ -gram overlap between the summary and the source. We demonstrated that it has high correlations with human judgement and can be used along other metrics to improve them, and to mitigate their sensitivity to low-quality references.

The prospects for future research include further exploration of the behaviour of reference-based, reference-free and hybrid metrics with references of varying quality, as well as potential extensions to multimodal settings such as the evaluation of vision-language systems.

## 7 Limitations

Like other lexical overlap metrics, ours works with the assumption that there is a vocabulary overlap between the source document and the summary, *ie* that the summary has a non-zero coverage. In order to evaluate the sensitivity of our metric to various levels of extractiveness of summaries, we would have wanted to compute the score on systems with varying values on the faithfulness-abstractiveness tradeoff curve presented in [Ladhak et al. \(2022\)](#); but their data was not made available yet.

[Vasilyev and Bohannon \(2021b\)](#) noticed that higher correlation with human scores can be achieved with "false" improvements, mimicking human behaviour. Using a referenceless evaluation metric, they limited the comparisons with the source text by selecting sentences to maximize their score, and observed a higher correlation with human judgement as a result. [Wu et al. \(2024\)](#) observe a similar consequence by first extracting sentences that maximize the ROUGE score with the original document and using the resulting extracted sentences along the predicted summary as the input to be evaluated by a LLM-as-a-judge. Their interpretation however is different as they do not view this higher correlation with human scores as a "false" improvement, but as a way to mitigate the *Lost-in-the-Middle* problem of LLMs.

We believe that the relevant interpretation depends on the method that is used to extract sentences from the source document. Using comparisons with the summary to extract "oracle" spans of the original document, or selecting key sentences that span over the main information of the document are not motivated by the same reasons. Mimicking the human behaviour of referring only to the bits of the document that are relevant to the proposed summary *at first glance* to score marginally higher correlations is a different thing than filtering the most important bits of a document relative to a measure of word importance.

Our metric filters out the  $n$ -grams with little semantic significance in the document. This can mimic the human bias of comparing the summary to salient sentences only, but it will also lower the influence of the artifacts of extractiveness discussed in section 3.

Our metric is also specific to the task of summarization and might correlate differently with human judgement on summarization tasks with different compression ratio, extractiveness, or style. Table 2

in the Appendix A illustrates this.

LLM-as-a-Judge methods can solve the issues of sensitivity to extractiveness and task settings, while providing more interpretable results, but are not exempt from biases and come with a noticeably higher cost.

## References

- Pierre Colombo, Maxime Peyrard, Nathan Noiry, Robert West, and Pablo Piantanida. 2023. [The Glass Ceiling of Automatic Evaluation in Natural Language Generation](#). In *Findings of the Association for Computational Linguistics: IJCNLP-AACL 2023 (Findings)*, pages 178–183, Nusa Dua, Bali. Association for Computational Linguistics.
- Daniel Deutsch, Rotem Dror, and Dan Roth. 2022. [Re-Examining System-Level Correlations of Automatic Summarization Evaluation Metrics](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 6038–6052, Seattle, United States. Association for Computational Linguistics.
- Esin Durmus, Faisal Ladhak, and Tatsunori Hashimoto. 2022. [Spurious Correlations in Reference-Free Evaluation of Text Generation](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1443–1454, 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. Place: Cambridge, MA Publisher: MIT Press.
- Manuel Faysse, Gautier Viaud, Céline Hudelot, and Pierre Colombo. 2023. [Revisiting Instruction Fine-tuned Model Evaluation to Guide Industrial Applications](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9033–9048, Singapore. Association for Computational Linguistics.
- Markus Freitag, David Grangier, and Isaac Caswell. 2020. [BLEU might be Guilty but References are not Innocent](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 61–71, Online. Association for Computational Linguistics.
- Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2023. [Prometheus: Inducing Fine-Grained Evaluation Capability in Language Models](#).

412	Seungone Kim, Juyoung Suk, Shayne Longpre,	Karen Sparck Jones. 1972. A statistical interpretation	470
413	Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham	of term specificity and its application in retrieval.	471
414	Neubig, Moontae Lee, Kyungjae Lee, and Minjoon	<i>Journal of documentation</i> , 28(1):11–21. Publisher:	472
415	Seo. 2024. <b>Prometheus 2: An Open Source Lan-</b>	MCB UP Ltd.	473
416	<b>guage Model Specialized in Evaluating Other Lan-</b>		
417	<b>guage Models.</b> <i>arXiv preprint</i> . ArXiv:2405.01535	Oleg Vasilyev and John Bohannon. 2021a. <b>ESTIME: Es-</b>	474
418	[cs].	<b>timation of Summary-to-Text Inconsistency by Mis-</b>	475
		<b>matched Embeddings.</b> In <i>Proceedings of the 2nd</i>	476
419	Huan Yee Koh, Jiaxin Ju, He Zhang, Ming Liu, and	<i>Workshop on Evaluation and Comparison of NLP</i>	477
420	Shirui Pan. 2022. <b>How Far are We from Robust</b>	<i>Systems</i> , pages 94–103, Punta Cana, Dominican Re-	478
421	<b>Long Abstractive Summarization?</b> In <i>Proceedings of</i>	public. Association for Computational Linguistics.	479
422	<i>the 2022 Conference on Empirical Methods in Nat-</i>		
423	<i>ural Language Processing</i> , pages 2682–2698, Abu	Oleg Vasilyev and John Bohannon. 2021b. <b>Is Human</b>	480
424	Dhabi, United Arab Emirates. Association for Com-	<b>Scoring the Best Criteria for Summary Evaluation?</b>	481
425	putational Linguistics.	In <i>Findings of the Association for Computational</i>	482
		<i>Linguistics: ACL-IJCNLP 2021</i> , pages 2184–2191,	483
426	Faisal Ladhak, Esin Durmus, He He, Claire Cardie,	Online. Association for Computational Linguistics.	484
427	and Kathleen McKeown. 2022. <b>Faithful or Extrac-</b>		
428	<b>tive? On Mitigating the Faithfulness-Abstractiveness</b>	Yunshu Wu, Hayate Iso, Pouya Pezeshkpour, Nikita	485
429	<b>Trade-off in Abstractive Summarization.</b> In <i>Proceed-</i>	Bhutani, and Estevam Hruschka. 2024. <b>Less is More</b>	486
430	<i>ings of the 60th Annual Meeting of the Association</i>	<b>for Long Document Summary Evaluation by LLMs.</b>	487
431	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	In <i>Proceedings of the 18th Conference of the Euro-</i>	488
432	<i>pers)</i> , pages 1410–1421, Dublin, Ireland. Association	<i>pean Chapter of the Association for Computational</i>	489
433	for Computational Linguistics.	<i>Linguistics (Volume 2: Short Papers)</i> , pages 330–343,	490
		St. Julian’s, Malta. Association for Computational	491
434	Chin-Yew Lin. 2004. <b>ROUGE: A Package for Auto-</b>	Linguistics.	492
435	<b>matic Evaluation of Summaries.</b> In <i>Text Summariza-</i>		
436	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.	Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021.	493
437	Association for Computational Linguistics.	<b>BARTScore: Evaluating Generated Text as Text Gen-</b>	494
		<b>eration.</b> In <i>Advances in Neural Information Process-</i>	495
438	Yixin Liu, Alex Fabbri, Pengfei Liu, Yilun Zhao, Liny-	<i>ing Systems</i> , volume 34, pages 27263–27277. Curran	496
439	ong Nan, Ruilin Han, Simeng Han, Shafiq Joty,	Associates, Inc.	497
440	Chien-Sheng Wu, Caiming Xiong, and Dragomir		
441	Radev. 2023. <b>Revisiting the Gold Standard: Ground-</b>	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.	498
442	<b>ing Summarization Evaluation with Robust Human</b>	Weinberger, and Yoav Artzi. 2019. <b>BERTScore:</b>	499
443	<b>Evaluation.</b> In <i>Proceedings of the 61st Annual Meet-</i>	<b>Evaluating Text Generation with BERT.</b>	500
444	<i>ing of the Association for Computational Linguistics</i>		
445	<i>(Volume 1: Long Papers)</i> , pages 4140–4170, Toronto,	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	501
446	Canada. Association for Computational Linguistics.	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	502
		Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang,	503
447	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and	Joseph E. Gonzalez, and Ion Stoica. 2023. <b>Judging</b>	504
448	Ryan McDonald. 2020. <b>On Faithfulness and Factu-</b>	<b>LLM-as-a-Judge with MT-Bench and Chatbot Arena.</b>	505
449	<b>ality in Abstractive Summarization.</b> In <i>Proceedings</i>		
450	<i>of the 58th Annual Meeting of the Association for</i>	Wanzheng Zhu and Suma Bhat. 2020. <b>GRUEN for Eval-</b>	506
451	<i>Computational Linguistics</i> , pages 1906–1919, On-	<b>uating Linguistic Quality of Generated Text.</b> In <i>Find-</i>	507
452	line. Association for Computational Linguistics.	<i>ings of the Association for Computational Linguistics:</i>	508
		<i>EMNLP 2020</i> , pages 94–108, Online. Association for	509
453	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Computational Linguistics.	510
454	Jing Zhu. 2002. <b>Bleu: a Method for Automatic Eval-</b>		
455	<b>uation of Machine Translation.</b> In <i>Proceedings of</i>	Maja Popović. 2015. <b>chrF: character n-gram F-score</b>	460
456	<i>the 40th Annual Meeting of the Association for Com-</i>	<b>for automatic MT evaluation.</b> In <i>Proceedings of the</i>	461
457	<i>putational Linguistics</i> , pages 311–318, Philadelphia,	<i>Tenth Workshop on Statistical Machine Translation</i> ,	462
458	Pennsylvania, USA. Association for Computational	pages 392–395, Lisbon, Portugal. Association for	463
459	Linguistics.	Computational Linguistics.	464
460	Maja Popović. 2015. <b>chrF: character n-gram F-score</b>		
461	<b>for automatic MT evaluation.</b> In <i>Proceedings of the</i>		
462	<i>Tenth Workshop on Statistical Machine Translation</i> ,		
463	pages 392–395, Lisbon, Portugal. Association for		
464	Computational Linguistics.		
465	S. E. Robertson and K. Sparck Jones. 1976.		
466	<b>Relevance weighting of search terms.</b> <i>Journal</i>		
467	<i>of the American Society for Informa-</i>		
468	<i>tion Science</i> , 27(3):129–146. <i>_eprint:</i>		
469	<a href="https://onlinelibrary.wiley.com/doi/pdf/10.1002/asi.4630270302">https://onlinelibrary.wiley.com/doi/pdf/10.1002/asi.4630270302</a> .		

## A Appendix

### A.1 Spurious correlations

Durmus et al. (2022) observed that model-based reference-free evaluation often have higher correlations with spurious correlates such as perplexity, length, coverage or density, than with human scores. We report the correlations between metrics and spurious correlates in Table 2.

### A.2 Correlations with human judgement on different settings

Figure 5 illustrate the distributions of system-level correlations of our metric with different settings.

For tokenization, we tested tokenizing texts as separated by space, using character tokenization, a pretrained GPT-2 tokenizer, or a custom tokenizer, trained on each corpus with a vocabulary of 100 tokens.

We included different sizes of  $n$ -grams in our tests, with bigrams, trigrams and 4-grams.

The two methods we considered for importance weighing are *tf-idf* and *bm-25*.

The importance score is the weight used to score the overlapped  $n$ -grams, we included the following scores:

- importance:  $t, d, D \mapsto w_{t,d,D}$
- exp-rank:  $t, d, D \mapsto \exp(-r_{t,d,D})$
- inv-rank:  $t, d, D \mapsto \frac{1}{r_{t,d,D}}$
- constant:  $t, d, D \mapsto 1$
- tanh:  $t, d, D \mapsto \tanh(\frac{w_{t,d,D}}{r_{t,d,D}})$

The options for the length penalty  $\alpha_{\hat{s},\hat{d}}$  are no penalty or  $\alpha_{\hat{s},d} = f(|\hat{s}|, |d|)$ , with

$$f : |\hat{s}|, |d| \mapsto \frac{1}{1 + \exp(20 * \frac{|\hat{s}|}{|d|} - 10)}$$

$f$  is illustrated in Figure 4.

We chose to use the corpus tokenizer, with trigrams, *tf-idf* and the tanh importance scoring with length penalty. These settings proved to be consistent in the tested conditions, and provided good ranges of values on different inputs. All the other experiments with our metric in this paper were using these settings.

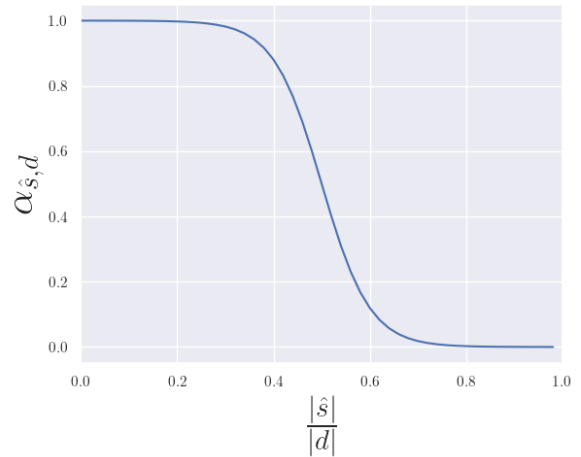


Figure 4: Length penalty  $\alpha_{\hat{s},d} = f(|\hat{s}|, |d|)$  with

$$f : |\hat{s}|, |\hat{d}| \mapsto \frac{1}{1 + \exp(20 * \frac{|\hat{s}|}{|\hat{d}|} - 10)}$$

### A.3 Range of values

We report the range of values taken by our metric, and ROUGE-1, for different inputs and on different datasets in Figures 6 and 7.

Table 2: Summary-level correlations between our metric, human evaluation metrics and spurious correlates. Values are bolded when the correlation with spurious correlate is higher than with human evaluation.

Metric	SummEval		arXiv		GovReport		RoSE	
	Pearson	Spearman	Pearson	Spearman	Pearson	Spearman	Pearson	Spearman
Relevance	0.24	0.23	0.42	0.45	0.15	0.18		
Coherence	0.20	0.23						
Consistency	0.03	0.04	-0.02	-0.11	-0.30	-0.32		
Fluency	0.01	-0.01						
Reference-based							0.17	0.14
Reference-free							0.18	0.15
Coverage	<b>0.26</b>	<b>0.28</b>	0.07	0.43	<b>-0.20</b>	<b>-0.19</b>	-0.08	0.05
Density	0.18	0.22	-0.04	-0.02	-0.03	0.17	0.04	0.02
Compression Ratio	-0.03	-0.06	0.01	0.02	<b>-0.54</b>	<b>-0.64</b>	<b>-0.28</b>	<b>-0.18</b>
Summary Length	<b>0.32</b>	<b>0.28</b>	0.40	0.28	0.02	-0.08	<b>0.30</b>	<b>0.26</b>



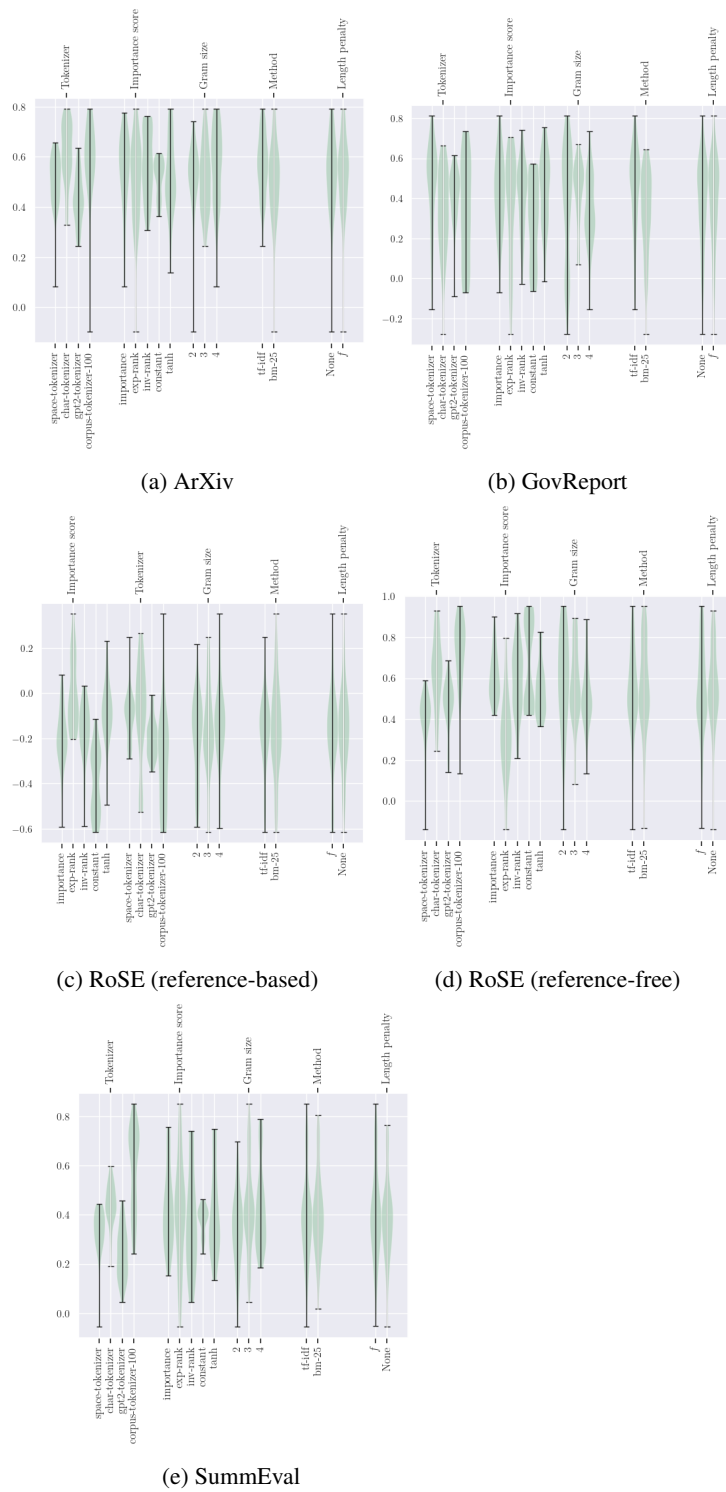
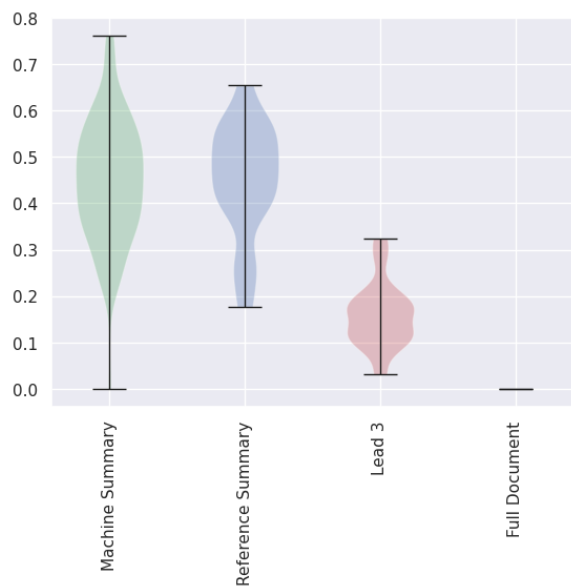
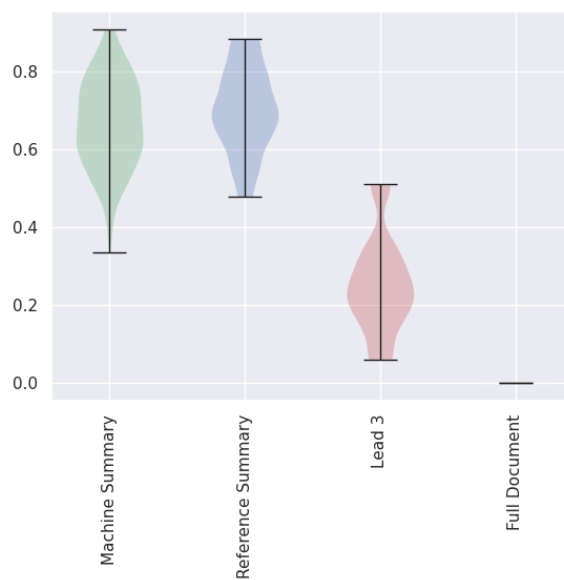


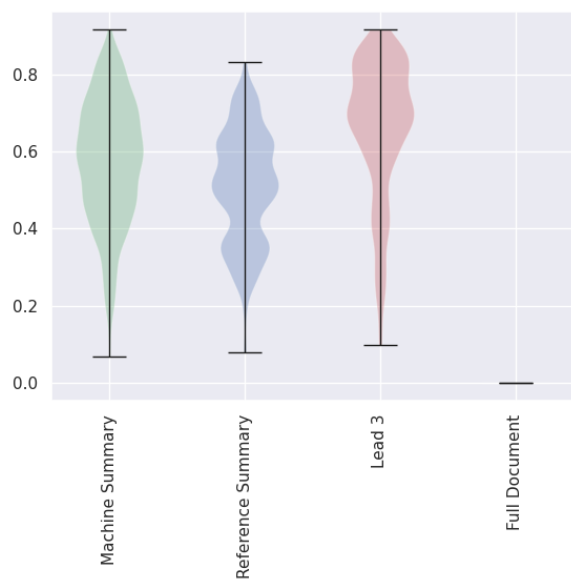
Figure 5: Distribution of system-level correlations of our metric in different settings



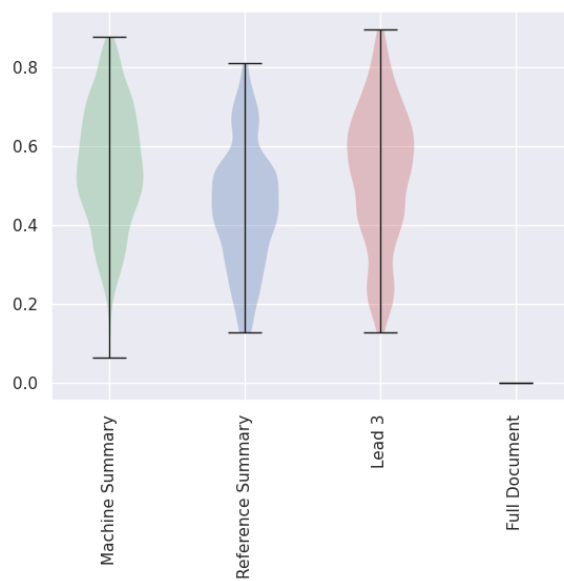
(a) ArXiv



(b) GovReport

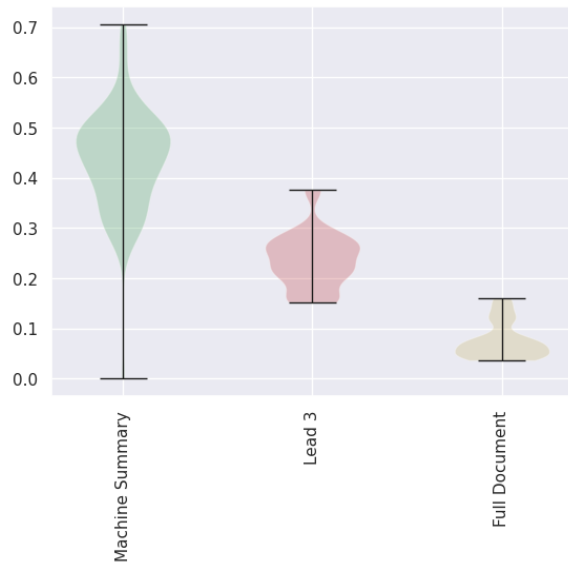


(c) SummEval

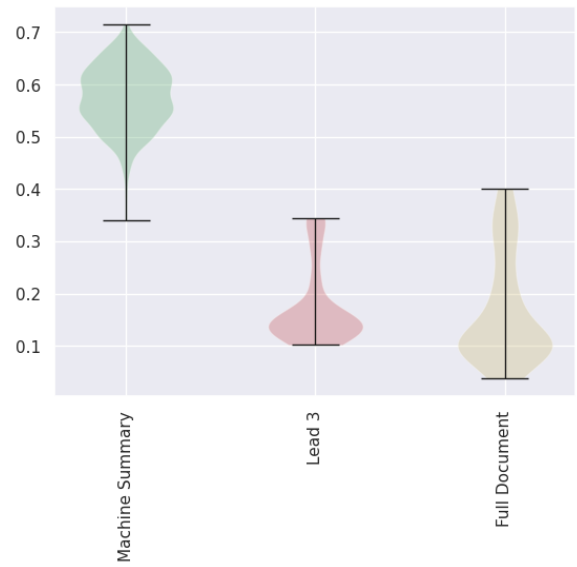


(d) RoSE

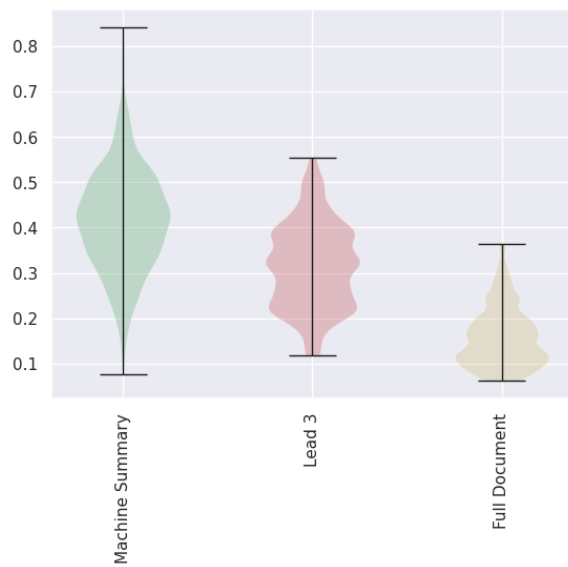
Figure 6: Range of values taken by our metric for different summaries



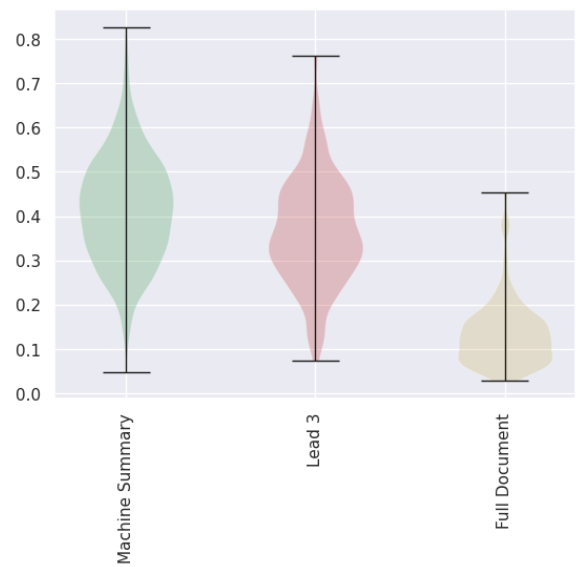
(a) ArXiv



(b) GovReport



(c) SummEval



(d) RoSE

Figure 7: Range of values taken by ROUGE-1 for different summaries