

# Stress Testing Factual Consistency Metrics for Long-Document Summarization

Anonymous ACL submission

## Abstract

Evaluating the factual consistency of abstractive text summarization remains a significant challenge, particularly for long documents, where conventional metrics struggle with input length limitations and long-range dependencies. In this work, we systematically evaluate the reliability of six widely used reference-free factuality metrics, originally proposed for short-form summarization, in the long-document setting. We probe metric robustness through seven factuality-preserving perturbations applied to summaries, namely paraphrasing, simplification, synonym replacement, logically equivalent negations, vocabulary reduction, compression, and source text insertion, and further analyze their sensitivity to retrieval context and claim information density. Across three long-form benchmark datasets spanning science fiction, legal, and scientific domains, our results reveal that existing short-form metrics produce inconsistent scores for semantically equivalent summaries and exhibit declining reliability for information-dense claims whose content is semantically similar to many parts of the source document. While expanding the retrieval context improves stability in some domains, no metric consistently maintains factual alignment under long-context conditions. Finally, our results highlight concrete directions for improving factuality evaluation, including multi-span reasoning, context-aware calibration, and training on meaning-preserving variations to enhance robustness in long-form summarization.<sup>1</sup>

## 1 Introduction

Abstractive summarization has seen rapid advances with the advent of large language models (LLMs), but ensuring that generated summaries faithfully reflect the source content remains a persistent challenge (Laban et al., 2024; Wright et al., 2025). Summaries that read fluently can

<sup>1</sup>We release our experimental code and all perturbed summary outputs to facilitate future research.

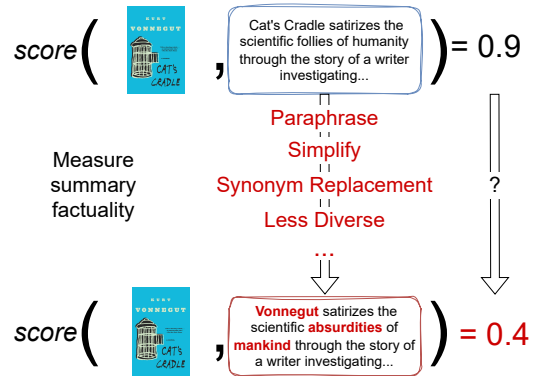


Figure 1: We aim to see how robust summary factuality metrics are for long and multi-document setups by applying meaning-preserving perturbations and comparing metric scores before and after these edits.

nonetheless introduce hallucinated details or omit critical facts (Belém et al., 2025), undermining their reliability for downstream tasks in domains such as medicine, law, and scientific review (Asai et al., 2024). Traditional evaluation measures like ROUGE (Lin, 2004) and BLEU (Papineni et al., 2002), which rely on n-gram overlap with reference summaries, are useful for measuring surface similarity but fail to capture factual consistency, since two summaries can overlap heavily in wording while still differ in correctness (Maynez et al., 2020). This gap has motivated the development of reference-free factuality evaluation metrics, which assess whether the statements in a summary are supported by the source document itself rather than by comparison with a human-written reference, using techniques such as question answering (Wang et al., 2020; Scialom et al., 2021; Fabbri et al., 2022), natural language inference (Laban et al., 2022; Chen and Eger, 2023; Zha et al., 2023), or LLM-based scoring (Liu et al., 2023; Fu et al., 2024).

While such metrics have demonstrated promise on short-document datasets, their scalability to

065	long-form summarization is likely to be hindered	integrate multi-span reasoning and context-aware	117
066	by challenges unique to long context lengths (Rus-	calibration.	118
067	sak et al., 2024; Sarthi et al., 2024; Edge et al.,		
068	2024). Important details may be dispersed across	<b>2 Factual Consistency in Abstractive</b>	119
069	hundreds or thousands of tokens and thus over-	<b>Summarization</b>	120
070	looked by metrics that process only truncated in-		
071	puts (Laban et al., 2024); multi-document sum-	Factual consistency refers to whether a summary	121
072	maries must reconcile diverse writing styles and po-	accurately reflects the content of the source doc-	122
073	tentially conflicting information (Asai et al., 2024);	ument. While abstractive summarization sys-	123
074	and the reference-free setting deprives evaluators	tems have become increasingly fluent, they often	124
075	of gold annotations, necessitating robust intrinsic	produce factually incorrect or hallucinated state-	125
076	evaluation protocols. Our goal is to systematically	ments (Huang et al., 2025). These hallucinations	126
077	benchmark widely used factual consistency metrics	can range from minor misstatements to major dis-	127
078	under these long-document conditions, in order to	tortions, particularly when the input is lengthy and	128
079	reveal their robustness and limitations.	semantically dense (Belém et al., 2025). A num-	129
080	Building on this foundation, we apply a	ber of techniques have been proposed to mitigate	130
081	stress-testing methodology established in previ-	this (Gao et al., 2023; Qiu et al., 2023; Zhang et al.,	131
082	ous work (Ramprasad and Wallace, 2024) to six	2024; Mündler et al., 2024; Wang et al., 2024). De-	132
083	widely-used reference-free metrics (§ 4) across	spite these efforts, most prior work has focused	133
084	seven factuality-preserving perturbations reflect-	on short-form summarization settings, where the	134
085	ing realistic long-form summarization phenomena.	task is relatively controlled. In contrast, long-form	135
086	These perturbations (§ 3.1) broadly cover para-	summarization requires condensing information	136
087	phrasing operations, simplification of complex con-	spread across thousands of tokens, making reliable	137
088	structions, synonym replacement, reduced lexical	factuality evaluation substantially more difficult,	138
089	diversity, logically equivalent negations, further	especially without reference summaries or human	139
090	compression of the summary, and insertion of un-	annotations.	140
091	related source sentences. They are designed to		
092	challenge metrics to distinguish genuine factual	<b>2.1 Factuality Evaluation Metrics</b>	141
093	consistency from superficial cues and to maintain		
094	robustness in the face of lexical and structural vari-	To overcome the limitations of traditional reference-	142
095	ation. On top of this, we investigate the impact	based metrics (Maynez et al., 2020), a range of	143
096	of long-document specific phenomena, including	reference-free metrics have emerged that assess	144
097	retrieval length and evidence dispersion (Goldman	the alignment between the summary and source	145
098	et al., 2024).	directly. Entailment-based approaches such as	146
099	Through extensive experiments on six evaluation	FactCC (Kryscinski et al., 2020) and SummaC (La-	147
100	metrics across three benchmark datasets in science	ban et al., 2022) use NLI models to judge whether	148
101	fiction, legal, and scientific domains, we expose	summary sentences are supported by the source.	149
102	significant weaknesses in existing metrics, such	QA-based methods like QAGS (Wang et al., 2020)	150
103	as inconsistent scoring across semantically equiva-	and QuestEval (Scialom et al., 2021) evaluate fac-	151
104	lency summaries (Fig. 1). Our analysis reveals that	tuality by generating and answering questions de-	152
105	current factuality metrics vary widely in robust-	derived from the summary. Generation-based met-	153
106	ness to meaning-preserving edits, with some highly	rics, including BARTScore (Yuan et al., 2021)	154
107	sensitive to surface changes while others remain	and T5Score (Trainin and Abend, 2025), esti-	155
108	more stable. Most metrics benefit from broader	mate the likelihood of the summary given the	156
109	retrieval context windows, though with notable	source using pretrained sequence-to-sequence mod-	157
110	domain-specific variation. We also find that metric	els. More recent tools like AlignScore (Zha	158
111	reliability decreases for information-dense claims	et al., 2023), MiniCheck (Tang et al., 2024), and	159
112	that overlap semantically with large portions of the	UniEval (Zhong et al., 2022) aim to improve ef-	160
113	source document, suggesting that current metrics	iciency and generalization across tasks. While	161
114	struggle with compressed or globally entangled	effective on short-document benchmarks, most of	162
115	content. These insights point toward improving	these metrics assume the full source and summary	163
116	evaluation consistency by developing metrics that	can be jointly encoded, limiting their utility for	164
		long-form inputs. Moreover, recent work shows	165

166	that these metrics are often brittle and sensitive	2.3 Adversarial Robustness of Metrics	216
167	to edits like paraphrasing, reordering, or logically	Recent work has evaluated the robustness of factu-	217
168	equivalent reformulations (Ramprasad and Wallace,	ality metrics by applying controlled perturbations	218
169	2024). A recent survey highlights persistent limita-	to the summary or source (Goyal and Durrett, 2021;	219
170	tions in robustness and long-document evaluation	Chen et al., 2021; Gabriel et al., 2021). Ramprasad	220
171	for factuality metrics (Lamsiyah et al., 2025). In	and Wallace (2024) showed that many metrics are	221
172	this work, we study the behavior of six popular	brittle when faced with logically equivalent but lex-	222
173	factuality metrics in long-document summariza-	ically altered summaries, with even benign transfor-	223
174	tion using a retrieval-based scoring framework, and	mations such as reordering or simplification caus-	224
175	systematically evaluate their robustness to a set of	ing large score shifts. However, these evaluations	225
176	controlled meaning-preserving perturbations.	have been primarily limited to short-document set-	226
177		ttings, where evidence is localized and both the	227
178	<b>2.2 Challenges in Long-Document Factuality</b>	source and summary can typically be processed	228
	<b>Evaluation</b>	jointly.	229
179	Evaluating factual consistency in long documents	Extending this evaluation framework to long-	230
180	introduces challenges that differ fundamentally	document summarization is not a straightforward	231
181	from those in short texts. Long inputs often contain	change of testbed. Long documents introduce ad-	232
182	information that is dispersed, hierarchically struc-	ditional challenges, including dispersed evidence,	233
183	tured, and cross-referential, requiring models to	higher degrees of abstraction and compression, and	234
184	link evidence across distant sections or even mul-	the need for retrieval-based evaluation to overcome	235
185	multiple documents (Asai et al., 2024). This results	input length constraints. These factors fundamen-	236
186	in long-range dependencies and positional biases	tally alter how factuality metrics operate and in-	237
187	such as the “lost in the middle” effect (Liu et al.,	teract with the input. In this work, we therefore	238
188	2024). Yet, research into robust factuality metrics	adopt the perturbation-based methodology of prior	239
189	for long inputs remains limited. Koh et al. (2023)	work as a foundation and systematically examine	240
190	identified a clear gap in the literature for automatic	how these metrics behave under long-context con-	241
191	evaluation methods tailored to long-document sum-	ditions. Beyond this, we additionally analyze the	242
192	marization. LongSciVerify (Bishop et al., 2024)	effects of retrieval context size and claim informa-	243
193	and LongEval (Krishna et al., 2023) are two of the	tion density, revealing new failure modes that arise	244
194	only available datasets with human factuality anno-	specifically in long-document and multi-document	245
195	tations in this setting. Chunk-based approaches like	settings. Our results show that brittleness observed	246
196	SMART (Amplayo et al., 2022), partially address	in short documents persists and is often amplified in	247
197	this by sequentially processing document segments,	long-form summarization, motivating the need for	248
198	but they remain computationally expensive and of-	factuality metrics that can reason over multi-span	249
199	ten inconsistent. A more scalable alternative is	evidence rather than relying on local or surface-	250
200	retrieval-based scoring, exemplified by LongDoc-	level cues.	251
201	FACTScore (Bishop et al., 2024), which retrieves	<b>3 Metric Robustness in Long-Form</b>	252
202	top-k relevant source passages for each summary	<b>Summarization</b>	253
203	sentence, and computes sentence-level factuality	To analyze the robustness of existing factuality met-	254
204	scores that can be aggregated into a global metric.	rics in long-form summarization, we evaluate six	255
205	However, it remains unclear whether existing met-	widely used reference-free metrics (§ 4) that span	256
206	rics, when used within such frameworks, behave	diverse architectures and scoring paradigms.	257
207	consistently and robustly across varying retrieval	<b>3.1 Perturbation Strategies</b>	258
208	configurations. Recent work suggests that uni-	To evaluate the robustness of factuality metrics in	259
209	formly aggregating sentence-level factuality scores	a controlled manner, we apply meaning-preserving	260
210	can be suboptimal for long documents, and that	perturbations to the original summaries, as done	261
211	discourse-aware aggregation can improve inconsis-	in Ramprasad and Wallace (2024) in the short doc-	262
212	tency detection (Zhong and Litman, 2025). Our	ument case. These perturbations are designed to	263
213	study addresses this gap by analyzing these behav-	vary the summary’s surface form (lexical choices,	264
214	iors under controlled conditions and across multi-		
215	ple domains.		

structure, or style) while preserving its factual consistency with the source document. In principle, a robust factuality metric should be invariant to such benign edits, assigning similar scores to the original and perturbed versions.

We define seven perturbation types, each targeting a different linguistic dimension. These include *Paraphrased*, where the summary is rewritten with alternate phrasings and syntactic structures; *Simplified*, where complex or compound constructions are rewritten into shorter, more readable sentences; *Synonym Replaced*, where content words are substituted with close synonyms to test for lexical invariance. We also generate *Less Diverse* summaries that reduce vocabulary variation, exploring whether metrics implicitly reward stylistic richness. Additional perturbations include *Negated*, which introduces logically equivalent negations to probe sensitivity to syntactic polarity, *Summarized*, which further compresses the summary to test how conciseness is handled, and *Added Source Text*, which inserts a factual sentence directly from the source that is unrelated to the main summary content. All seven perturbed summaries are generated using the GPT-4o (Hurst et al., 2024) model via the OpenAI API. The detailed prompts used to generate each perturbation are provided in App. A. To ensure that these perturbations preserve factual consistency, we additionally perform an NLI-based faithfulness check comparing each perturbed summary against its original counterpart; detailed results are reported in App. C.

While these transformations are meaning-preserving, they pose particular challenges in the long-document setting: summaries must capture information scattered across thousands of tokens, so perturbations that change sentence structure, reduce vocabulary, or alter flow can disrupt long-range dependencies and retrieval alignment. This makes them a rigorous test of whether factuality metrics remain robust. Any significant fluctuation in scores, despite no factual errors being introduced, indicates that a metric is reacting to surface-level edits rather than faithfully assessing factual consistency.

### 3.2 Retrieval-Based Scoring for Long Documents

Most factuality metrics in existing literature are designed for short inputs and cannot directly process the full content of long documents due to token length limitations. This is particularly problem-

atic in long-form summarization, where summaries may draw on information scattered across multiple sections or even multiple documents. To address this, we follow the retrieval-augmented strategy proposed by Bishop et al. (2024), which enables factuality evaluation at the sentence level without requiring the metric to ingest the entire source document at once.

Let  $S = \{s_1, s_2, \dots, s_m\}$  denote the summary, where each  $s_j$  is a sentence, and let  $D = \{d_1, d_2, \dots, d_n\}$  be the set of sentences in the source document. For each summary sentence  $s_j$ , we compute a sentence embedding  $e_j$ , and similarly obtain embeddings  $\{e_1^D, \dots, e_n^D\}$  for the source document using a pre-trained sentence encoder. We compute cosine similarity<sup>2</sup> between  $s_j$  and each  $d_i \in D$  and retrieve the top- $K$  most similar source sentences. Each retrieved sentence  $d_{j,k}$  is then expanded to include the surrounding context within a symmetric window size  $w$ , forming a snippet:

$$d_{j,k}^{(w)} = \{d_{j,k-w}, \dots, d_{j,k}, \dots, d_{j,k+w}\}. \quad (1)$$

We evaluate the factual consistency of  $s_j$  against each of the  $K$  context snippets, using any automated metric  $\mathcal{M}$ , and take the maximum score:

$$\text{score}(s_j) = \max_{k \in \{1, \dots, K\}} \mathcal{M}(s_j, d_{j,k}^{(w)}). \quad (2)$$

Finally, the summary-level factuality score is computed by averaging these sentence-level scores across all sentences in the summary. In our experiments, we explore how varying  $w$  affects metric behavior, shedding light on the sensitivity of different metrics to retrieval context size.

## 4 Experimental Setup

**Metrics** We evaluate six reference-free factuality metrics that represent diverse architectures and scoring paradigms. BARTScore (Yuan et al., 2021) estimates the log-likelihood of the summary given the source using a pretrained BART model<sup>3</sup>, treating factuality as a conditional generation problem. For consistency with other metrics, we exponentiate the log-likelihood scores to obtain normalized values that are directly comparable across metrics. SummaC-Conv and SummaC-ZS (Laban et al., 2022)

<sup>2</sup>We use SBERT (Reimers and Gurevych, 2019) for computing all sentence similarities.

<sup>3</sup><https://huggingface.co/facebook/bart-large-cnn>

represent entailment-based approaches that apply pretrained NLI models to compute consistency between summary and source sentence pairs; the former uses a learned aggregation layer, while the latter relies on zero-shot averaging. *AlignScore* (Zha et al., 2023) leverages contrastive alignment learning across NLI, QA, and summarization tasks and demonstrates strong cross-domain generalization. *UniEval* (Zhong et al., 2022) formulates the evaluation as a multi-dimensional question answering task within a unified T5-based framework, jointly considering factuality, coherence, relevance, and fluency. Finally, *MiniCheck* (Tang et al., 2024) is a lightweight, sentence-level factuality classifier that achieves near GPT-4 performance at a fraction of the cost. We use the *Bespoke-MiniCheck-7B* variant, which ranks highest on the *LLM-AggreFact* benchmark (Tang et al., 2024), and take advantage of its 32k-token context window to provide full-document context during evaluation. We evaluate these metrics in their publicly released form, reflecting common practice in prior studies that apply summarization metrics without task-specific adaptation (Huang et al., 2021; Yang and Wan, 2022; Sotudeh et al., 2021; Guo et al., 2022).

**Datasets** We conduct our analysis on three long-document summarization datasets spanning diverse domains: *SQuALITY* (Wang et al., 2022), *LexAbSumm* (T.y.s.s. et al., 2024), and *ScholarQABench* (Asai et al., 2024). *SQuALITY* consists of public domain science fiction stories paired with expert-written summaries that balance narrative abstraction and fine-grained detail. *LexAbSumm* contains legal judgments from the European Court of Human Rights, where summaries are aspect-specific and demand precise distillation of dense legal arguments. *ScholarQABench* is a multi-document benchmark based on open-access computer science papers, where the task is to generate detailed, factual answers to expert-written queries using evidence from multiple documents. We selected these three distinct datasets to ensure broad cross-domain coverage, given their substantial differences in structure, style, and language. Detailed dataset statistics are provided in App. B.

**Experimental Design** To evaluate metric robustness, we construct perturbed versions of the summaries from each dataset using the seven meaning-preserving transformations described in § 3.1. These edits allow us to test whether metrics exhibit sensitivity to benign changes. For each sum-

mary, original and perturbed, we use the framework (§ 3.2) to retrieve source evidence and compute factuality scores using all six metrics discussed above. Beyond this robustness analysis, we also investigate how retrieval granularity and claim information density influence metric behavior. We vary the evidence window size  $w$  from Eq. 1 to test how broader or narrower retrieval contexts affect metric performance. Additionally, we measure the information density of each summary sentence using the mean pairwise cosine similarity between its embedding and all sentences in the source document. High information density indicates claims that semantically overlap with many parts of the source and are therefore harder to verify, while low-density claims correspond to specific statements with localized evidence. This analysis reveals how metrics respond to different levels of semantic compression, providing insight into their sensitivity to claim complexity in long-form summarization.

## 5 Results & Analysis

We evaluate the robustness and behavior of six reference-free factuality metrics across a range of semantic-preserving perturbations, varying retrieval context windows, and differences in claim information density. Our results are presented in three parts: (1) robustness under perturbation (§ 5.1), (2) sensitivity to retrieval granularity (§ 5.2), and (3) metric sensitivity to claim information density (§ 5.3).

### 5.1 Robustness Against Perturbations

Fig. 2 presents the change in factuality scores when different perturbations are applied to the summaries across three datasets. Each plot shows the difference in score (perturbed minus original) for a given metric, broken down by perturbation and dataset. A robust metric should remain invariant to these semantic-preserving changes. However, we find that all metrics show varying levels of sensitivity.

On *LexAbSumm*, *BARTScore* shows clear negative shifts across nearly all perturbations, while remaining relatively consistent on all other datasets. These consistent declines on *LexAbSumm* indicate that *BARTScore* is highly sensitive to even mild surface-level edits in the legal domain, where long and complex sentence structures and domain-specific jargon likely amplify generation-based instability. *MiniCheck* shows very small changes across all perturbations and datasets. However, it

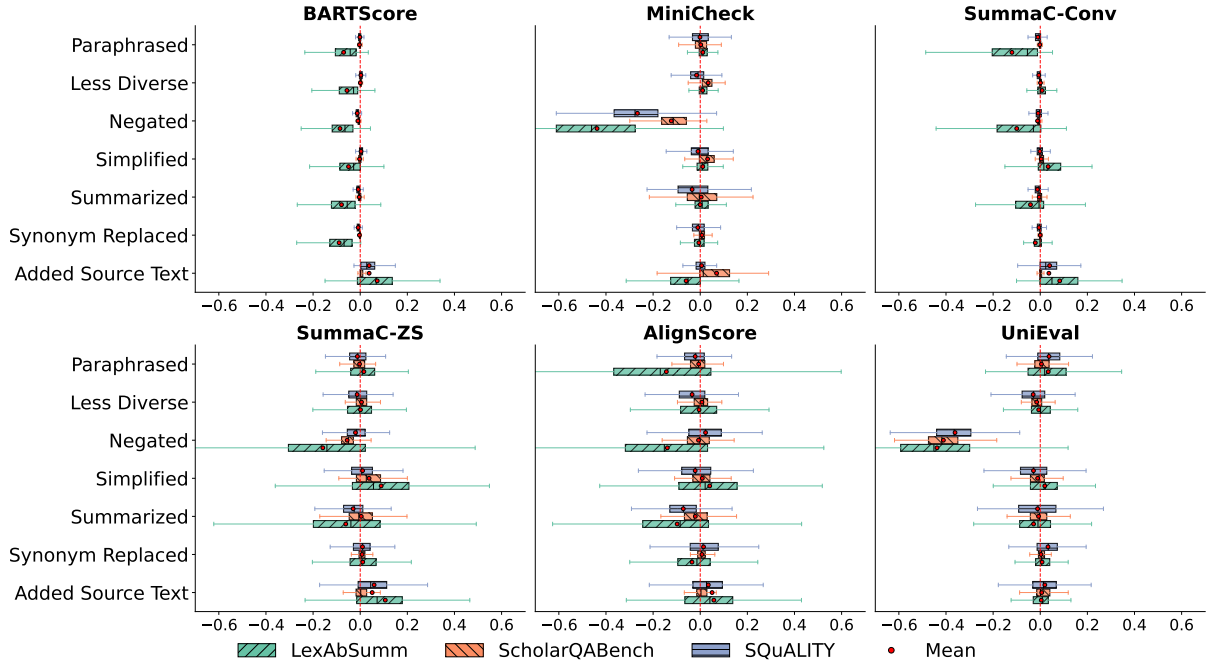


Figure 2: Score change under factuality-preserving perturbations. Boxplots show the difference in factuality score between the perturbed and original summaries, for each metric and perturbation type, across three datasets. The central dot indicates the mean score difference, and the whiskers represent the minimum and maximum values.

struggles with logically equivalent negations, especially in *LexAbSumm*. This may reflect a domain mismatch, as the metric appears less effective in capturing factual consistency in legal texts. SummaC-Conv and SummaC-ZS, both based on NLI, show moderate and more balanced behavior. They are somewhat affected by *Summarized* and *Negated* summaries, especially on *SQuALITY* and *LexAbSumm*, showing they are not fully invariant to meaning-preserving rewrites. UniEval is sensitive to most of the perturbations. However, it consistently fails to handle logically equivalent *Negated* summaries across all datasets, suggesting a lack of sensitivity to logical form. AlignScore mostly struggles with the legal domain, showing large score drops in response to *Paraphrased*, *Negated* and *Summarized* summaries. This suggests difficulty in tracking sentence order and logical consistency in structured, formal texts. While it performs more reliably on *SQuALITY* and *ScholarQABench*, it remains less robust than other metrics overall.

Detailed per-dataset results are provided in App. D. These results list the mean factuality scores for each metric and perturbation type across all datasets. To quantify domain-specific instability that is masked by signed averages, we also analyze mean absolute score changes under perturbations; detailed results are provided in App. E.

## 5.2 Effect of Retrieval Context Window Size

Table 1 reports average factuality scores on the original summaries for each metric and dataset, using window sizes  $w = 0, 1, 2$  in Eq. 1. We find that most metrics show consistent improvements as the window size increases, suggesting that they can effectively leverage broader local context when making sentence-level factuality judgments. This pattern is particularly pronounced on *LexAbSumm*, where understanding legal arguments often requires attending to multi-sentence spans. SummaC-ZS and SummaC-Conv display little sensitivity to larger context windows, implying that their underlying models base judgments on more localized comparisons and are less responsive to extended evidence.

Taken together, these results indicate that retrieval-based scoring can improve factuality assessment in long-document summarization, especially when a broader context is provided. However, NLI-based metrics remain insensitive to increasing context windows.

## 5.3 Metric Sensitivity to Claim Similarity

To understand what makes factual consistency evaluation difficult in long documents, we analyze how the semantic density of a claim relates to metric reliability. We hypothesize that claims which are highly similar to many claims in the original docu-

Metric	ScholarQABench			SQuALITY			LexAbSumm		
	$w = 0$	$w = 1$	$w = 2$	$w = 0$	$w = 1$	$w = 2$	$w = 0$	$w = 1$	$w = 2$
BARTScore	0.03	0.03	0.02	0.03	0.03	0.03	0.15	0.16	0.16
MiniCheck	0.17	0.15	0.15	0.11	0.15	0.19	0.47	0.53	0.60
SummaC-Conv	0.22	0.23	0.25	0.22	0.23	0.24	0.33	0.33	0.34
SummaC-ZS	0.14	0.18	0.20	0.11	0.13	0.14	0.36	0.38	0.39
AlignScore	0.15	0.21	0.27	0.10	0.18	0.24	0.36	0.52	0.64
UniEval	0.72	0.73	0.74	0.67	0.68	0.70	0.81	0.82	0.84

Table 1: Impact of retrieval context window size  $w$  on factuality scores for original summaries. Increasing  $w$  expands the snippet around each retrieved sentence.

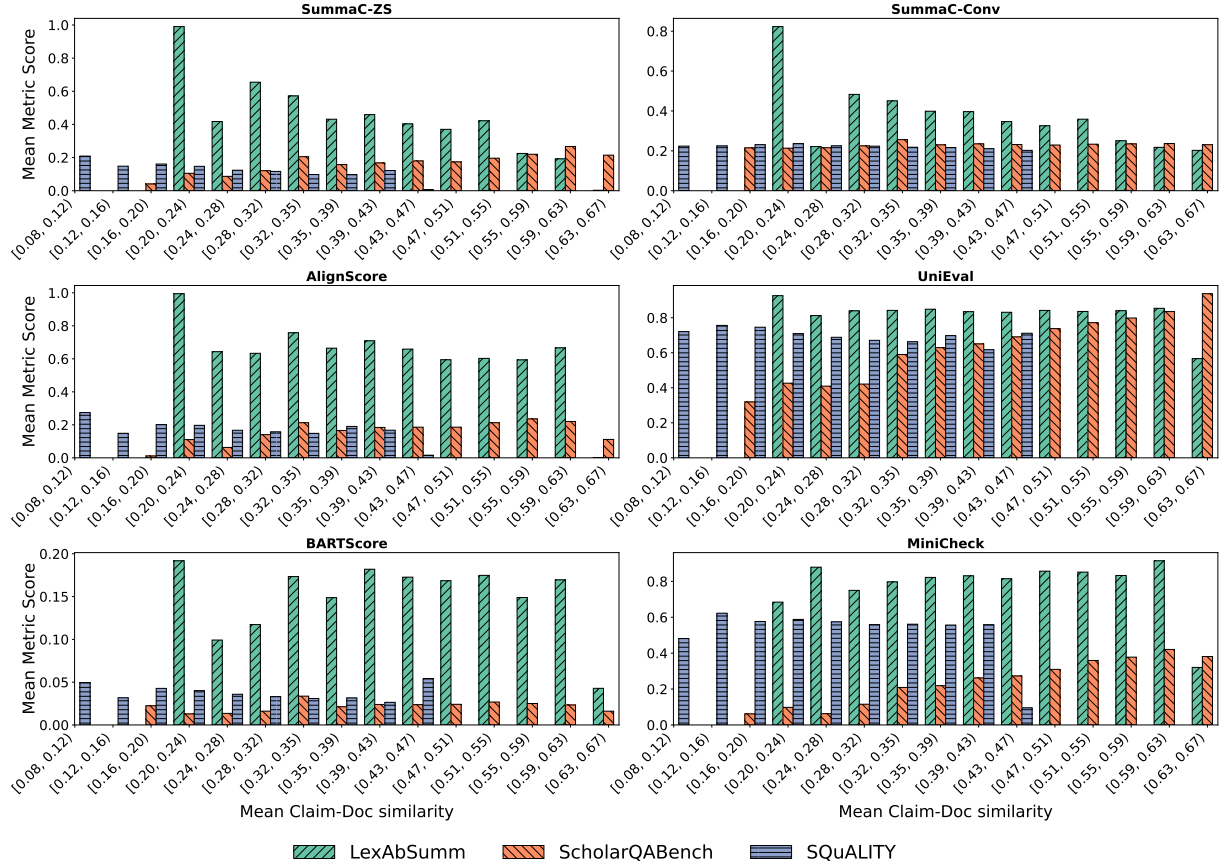


Figure 3: Relationship between claim similarity and average factuality score. Higher similarity values correspond to more information-dense claims whose content overlaps with multiple parts of the source document. Metrics generally assign lower scores to these claims for LexAbSumm and SQuALITY, and higher scores for ScholarQABench, indicating reduced reliability for compressed information.

ment contain more general (and thus compressed) information, and so are harder to fact-check since evidence for them is dispersed across multiple parts of the source document (Goldman et al., 2024).

We approximate this by computing the mean pairwise cosine similarity between each summary sentence  $s_j$  (claim) and all  $n$  sentences  $d_i$  in the source document  $D$ ,

$$\text{Sim}(s_j, D) = \frac{1}{n} \sum_{i=1}^n \cos(e_j, e_i^D), \quad (3)$$

where  $e_j$  and  $e_i^D$  denote the sentence embeddings of the claim and document sentences, respectively. Higher  $\text{Sim}(s_j, D)$  values indicate more broad claims whose meaning overlaps widely with the source, while lower values correspond to specific, localized claims. Claims are then grouped into similarity bins  $\mathcal{B}$ , and the average factuality score for each bin is calculated as

$$\text{Score}_{\text{bin}} = \frac{1}{|\mathcal{B}|} \sum_{s_j \in \mathcal{B}} M(s_j), \quad (4)$$

where  $M(s_j)$  is the factuality score for claim  $s_j$ .

A few trends emerge between claim similarity and metric sensitivity across all datasets and metrics, as shown in Fig. 3. For *LexAbSumm*, metric scores consistently decrease as claim similarity increases, meaning that the more a claim’s meaning is entangled with the broader document, the worse the metric becomes at predicting factuality. This is likely because summaries of legal documents may refer to specific aspects of the texts which are easier to fact-check, while more general statements which compress a lot of technical language are more difficult. The same occurs (though less pronounced) for *SQuALITY*, where more general claims are more challenging to fact check. In this case, we are summarizing novels, so more general claims try to compress the story narrative into a compact form, and thus will require retrieving or attending to and reasoning over disparate pieces of the text.

This effect is particularly visible in *AlignScore* and *BARTScore*, both of which depend on local lexical alignment or sentence-level contextual matching. Their scores show sharp declines for high-similarity claims, reflecting their sensitivity to distributed evidence. *SummaC-Conv* and *SummaC-ZS*, while somewhat more stable, also exhibit a gradual drop as similarity increases, which suggests that NLI-based judgments still rely on relatively localized entailment cues. In contrast, *UniEval* and *MiniCheck* maintain comparatively stable performance across bins, implying a higher degree of robustness to distributed or compressed content, although some degradation is still observed in the highest-similarity regions.

On the contrary, we see that more general statements are easier to fact check in *ScholarQABench* for many metrics. This could be because *ScholarQABench* is multi-document, so many sentences being similar to one claim may actually simply be repeated instances of the claim across documents. We see this upward trend especially on *UniEval* and *MiniCheck*, where metric quality is highly dependent on how general or specific a claim is.

On the whole, these findings support the broader conclusion that factuality metrics are dependent on how dispersed and overlapping the evidence is for a given claim, a hallmark of long-document summarization. Improving robustness for such cases likely requires metrics that can reason over multi-span evidence rather than relying on local or pairwise semantic alignment.

## 6 Conclusion & Future Work

In this work, we present a comprehensive evaluation of six widely used reference-free factuality metrics: *BARTScore*, *SummaC-Conv*, *SummaC-ZS*, *AlignScore*, *MiniCheck*, and *UniEval*. We tested their behavior under seven meaning-preserving perturbations applied to long-document summaries to assess whether these metrics reliably capture factual consistency. To enable evaluation over long documents, we used a sentence-level retrieval-based scoring strategy, which compares each summary sentence to the most relevant evidence snippets from the source document. This setup enabled fine-grained evaluation across three diverse long-form abstractive summarization datasets: *SQuALITY*, *LexAbSumm*, and *ScholarQABench*, covering sci-fi, legal, and scientific domains.

Our results revealed that many metrics respond inconsistently to perturbations that do not affect factual consistency. Several metrics exhibit unstable behavior in response to paraphrasing, simplification, and logically equivalent negations. *AlignScore* and *SummaC-ZS* are particularly unreliable across domains and perturbation types. In contrast, *UniEval* and *MiniCheck* are relatively robust, although they too struggle in specific cases, such as handling logical negations. Most metrics improve when evidence retrieval windows are expanded, particularly for complex, multi-sentence inputs such as legal documents. We also found that metrics are systematically affected by the information-density of claims whose meaning overlaps broadly with the source document, indicating that current approaches struggle to evaluate compressed or contextually entangled statements, which are common in long-form summarization. These findings highlight a need for factuality metrics that are robust to stylistic and logical variation, retrieval-aware, and sensitive to information density.

Future work should explore multi-span reasoning and context-aware calibration to better model distributed evidence, as well as contrastive training on meaning-preserving perturbations to improve stability. Incorporating human judgments can identify systematic weaknesses and support the design of metrics that generalize across domains and languages. We also see potential in hybrid approaches that combine reference-free with reference-based alignment signals, bridging semantic precision with contextual coverage for more reliable evaluation of long-document summarization.

## 635 Limitations

636 While our study offers a systematic investiga-  
637 tion into the robustness of factuality metrics un-  
638 der meaning-preserving perturbations in long-  
639 document summarization, there are several aspects  
640 that merit further consideration. Our analysis relies  
641 on automatically generated perturbations, produced  
642 using GPT-4o, which are designed to preserve fac-  
643 tual consistency. However, without human annota-  
644 tions, we cannot confirm with full certainty that all  
645 edits preserve factual correctness in every case. We  
646 also do not evaluate metric outputs against human  
647 factuality judgments in the long-document setting.  
648 Large-scale human annotations for long-document  
649 summaries are currently scarce, and conducting  
650 such an evaluation would require the creation of  
651 a new benchmark with human judgments across  
652 multiple metrics, domains, and perturbation types.  
653 This represents a substantial research effort beyond  
654 the scope of this work. This may introduce noise  
655 into the interpretation of metric behavior, especially  
656 when changes are subtle or domain-specific. We  
657 evaluate six reference-free metrics in their origi-  
658 nal, publicly released form and do not investigate  
659 whether fine-tuning, calibration, or adaptation to  
660 long-form inputs might mitigate some of the ob-  
661 served weaknesses. In our retrieval-based scoring  
662 setup, we use a fixed number of top-k most simi-  
663 lar sentences and vary only the surrounding con-  
664 text window to control retrieval granularity. While  
665 this gives us insight into how context affects met-  
666 ric behavior, it assumes a static retrieval strategy  
667 and does not account for dynamic query-based re-  
668 trieval or more sophisticated evidence selection  
669 methods that may better match human annotation  
670 patterns. Additionally, our analysis is confined  
671 to English-language datasets from three domains:  
672 science fiction, legal text, and scientific articles.  
673 These domains offer diversity in structure and style,  
674 but our findings may not fully generalize to other  
675 high-stakes applications such as medical or finan-  
676 cial summarization, or to non-English and low-  
677 resource settings. Addressing these broader limi-  
678 tations will be important for future work aiming  
679 to build more generalizable and reliable factuality  
680 evaluation pipelines.

## 681 Ethical Implications

682 This study evaluates automatic factuality metrics  
683 rather than developing new summarization models,  
684 and thus presents minimal direct ethical risk. How-

ever, factuality evaluation plays an important role  
in ensuring the reliability of language model out-  
puts. Weak or biased metrics could inadvertently  
overestimate the truthfulness of generated content,  
particularly in sensitive domains such as medicine  
or law. By identifying systematic weaknesses and  
proposing strategies for more reliable evaluation,  
this work aims to support safer and more account-  
able deployment of summarization systems. All  
datasets used in this study are publicly available  
and contain no personally identifiable information.

## References

- Reinold Kim Amplayo, Peter J Liu, Yao Zhao, and  
Shashi Narayan. 2022. [SMART: Sentences as  
basic units for text evaluation](#). *arXiv preprint  
arXiv:2208.01030*.
- Akari Asai, Jacqueline He, Rulin Shao, Weijia Shi,  
Amanpreet Singh, Joseph Chee Chang, Kyle Lo,  
Luca Soldaini, Sergey Feldman, Mike D’Arcy,  
David Wadden, Matt Latzke, Minyang Tian, Pan Ji,  
Shengyan Liu, Hao Tong, Bohao Wu, Yanyu Xiong,  
Luke Zettlemoyer, and 6 others. 2024. [OpenScholar:  
Synthesizing Scientific Literature with Retrieval-  
augmented LMs](#). *arXiv preprint arXiv:2411.14199*.
- Catarina G. Belém, Pouya Pezeshkpour, Hayate Iso,  
Seiji Maekawa, Nikita Bhutani, and Estevam Hr-  
uschka. 2025. [From single to Multi: How llms hal-  
lucinate in multi-document summarization](#). In *Find-  
ings of the Association for Computational Linguistics:  
NAACL 2025, Albuquerque, New Mexico, USA, April  
29 - May 4, 2025*, pages 5276–5309. Association for  
Computational Linguistics.
- Steven Bird. 2006. [NLTK: the natural language toolkit](#).  
In *ACL 2006, 21st International Conference on Com-  
putational Linguistics and 44th Annual Meeting of  
the Association for Computational Linguistics, Pro-  
ceedings of the Conference, Sydney, Australia, 17-21  
July 2006*. The Association for Computer Linguistics.
- Jennifer Bishop, Sophia Ananiadou, and Qianqian Xie.  
2024. [LongDocFACTScore: Evaluating the factu-  
ality of long document abstractive summarisation](#).  
In *Proceedings of the 2024 Joint International Con-  
ference on Computational Linguistics, Language Re-  
sources and Evaluation, LREC/COLING 2024, 20-25  
May, 2024, Torino, Italy*, pages 10777–10789. ELRA  
and ICCL.
- Yanran Chen and Steffen Eger. 2023. [MENLI: Robust  
evaluation metrics from natural language inference](#).  
*Transactions of the Association for Computational  
Linguistics*, 11:804–825.
- Yiran Chen, Pengfei Liu, and Xipeng Qiu. 2021. [Are factuality checkers reliable? adversarial meta-  
evaluation of factuality in summarization](#). In *Find-  
ings of the Association for Computational Linguis-  
tics: EMNLP 2021*, pages 2082–2095, Punta Cana,





968	Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020.	Shuo Zhang, Liangming Pan, Junzhou Zhao, and	1024
969	<a href="#">Asking and answering questions to evaluate the fac-</a>	William Yang Wang. 2024. <a href="#">The knowledge align-</a>	1025
970	<a href="#">tual consistency of summaries</a> . In <i>Proceedings of the</i>	<a href="#">ment problem: Bridging human and external knowl-</a>	1026
971	<i>58th Annual Meeting of the Association for Compu-</i>	<a href="#">edge for large language models</a> . In <i>Findings of the As-</i>	1027
972	<i>tational Linguistics</i> , pages 5008–5020, Online. Asso-	<i>sociation for Computational Linguistics, ACL 2024,</i>	1028
973	ciation for Computational Linguistics.	<i>Bangkok, Thailand and virtual meeting, August 11-</i>	1029
974	Alex Wang, Richard Yuanzhe Pang, Angelica Chen, Ja-	<i>16, 2024</i> , pages 2025–2038. Association for Compu-	1030
975	son Phang, and Samuel R. Bowman. 2022. <a href="#">SQuAL-</a>	tational Linguistics.	1031
976	<a href="#">ITY: Building a long-document summarization</a>	Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu	1032
977	<a href="#">dataset the hard way</a> . In <i>Proceedings of the 2022 Con-</i>	Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and	1033
978	<i>ference on Empirical Methods in Natural Language</i>	Jiawei Han. 2022. <a href="#">Towards a unified multi-</a>	1034
979	<i>Processing</i> , pages 1139–1156, Abu Dhabi, United	<a href="#">dimensional evaluator for text generation</a> . In <i>Pro-</i>	1035
980	Arab Emirates. Association for Computational Lin-	<i>ceedings of the 2022 Conference on Empirical Meth-</i>	1036
981	guistics.	<i>ods in Natural Language Processing</i> , pages 2023–	1037
982	Yuxia Wang, Revanth Gangi Reddy, Zain Muhammad	2038, Abu Dhabi, United Arab Emirates. Association	1038
983	Mujahid, Arnav Arora, Aleksandr Rubashevskii, Ji-	for Computational Linguistics.	1039
984	ahui Geng, Osama Mohammed Afzal, Liangming	Yang Zhong and Diane J. Litman. 2025. <a href="#">Discourse-</a>	1040
985	Pan, Nadav Borenstein, Aditya Pillai, Isabelle Au-	<a href="#">driven evaluation: Unveiling factual inconsistency</a>	1041
986	genstein, Iryna Gurevych, and Preslav Nakov. 2024.	<a href="#">in long document summarization</a> . In <i>Proceedings</i>	1042
987	<a href="#">Factcheck-bench: Fine-grained evaluation bench-</a>	<i>of the 2025 Conference of the Nations of the Amer-</i>	1043
988	<a href="#">mark for automatic fact-checkers</a> . In <i>Findings of the</i>	<i>icas Chapter of the Association for Computational</i>	1044
989	<i>Association for Computational Linguistics: EMNLP</i>	<i>Linguistics: Human Language Technologies, NAACL</i>	1045
990	2024, pages 14199–14230, Miami, Florida, USA.	<i>2025 - Volume 1: Long Papers, Albuquerque, New</i>	1046
991	Association for Computational Linguistics.	<i>Mexico, USA, April 29 - May 4, 2025</i> , pages 2050–	1047
992	Dustin Wright, Zain Muhammad Mujahid, Lu Wang,	2073. Association for Computational Linguistics.	1048
993	Isabelle Augenstein, and David Jurgens. 2025. Un-		
994	structured Evidence Attribution for Long Context		
995	Query Focused Summarization. In <i>Proceedings of</i>		
996	<i>the 2025 Conference on Empirical Methods in Natu-</i>		
997	<i>ral Language Processing (EMNLP)</i> . Association for		
998	Computational Linguistics.		
999	Cai Yang and Stephen Wan. 2022. <a href="#">Investigating metric</a>		
1000	<a href="#">diversity for evaluating long document summarisa-</a>		
1001	<a href="#">tion</a> . In <i>Proceedings of the Third Workshop on Schol-</i>		
1002	<i>arly Document Processing, SDP@COLING 2022,</i>		
1003	<i>Gyeongju, Republic of Korea, October 12 - 17, 2022,</i>		
1004	pages 115–125. Association for Computational Lin-		
1005	guistics.		
1006	Joonho Yang, Seunghyun Yoon, ByeongJeong Kim, and		
1007	Hwanhee Lee. 2024. <a href="#">FIZZ: Factual inconsistency</a>		
1008	<a href="#">detection by zoom-in summary and zoom-out docu-</a>		
1009	<a href="#">ment</a> . In <i>Proceedings of the 2024 Conference on</i>		
1010	<i>Empirical Methods in Natural Language Processing,</i>		
1011	pages 30–45, Miami, Florida, USA. Association for		
1012	Computational Linguistics.		
1013	Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021.		
1014	<a href="#">BARTScore: Evaluating Generated Text as Text Gen-</a>		
1015	<a href="#">eration</a> . In <i>Advances in Neural Information Process-</i>		
1016	<i>ing Systems (NeurIPS)</i> .		
1017	Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu.		
1018	2023. <a href="#">AlignScore: Evaluating factual consistency</a>		
1019	<a href="#">with a unified alignment function</a> . In <i>Proceedings</i>		
1020	<i>of the 61st Annual Meeting of the Association for</i>		
1021	<i>Computational Linguistics (Volume 1: Long Papers),</i>		
1022	pages 11328–11348, Toronto, Canada. Association		
1023	for Computational Linguistics.		

## A List of Prompts

The following prompts in Figure 4 are used with GPT-4o to produce each of the seven meaning-preserving perturbations described in § 3.1. Each prompt instructs the model to rewrite the summary according to a specific linguistic transformation while preserving factual meaning.

## B Dataset Statistics

The detailed statistics for the datasets used in our experiments can be seen in Table 2.

## C Faithfulness of Perturbed Summaries

To verify that the applied perturbations preserve factual consistency, we perform an automatic faithfulness check using an NLI-based approach. This analysis serves as a sanity check to ensure that the perturbations do not introduce widespread factual errors, rather than as a definitive evaluation of summary correctness.

For each perturbed summary, we split its text into sentences, evaluate it against the corresponding original summary, treating the original summary as the premise and the perturbed sentence as the hypothesis. If a premise–hypothesis pair exceeds the model’s maximum input length, we apply sentence-level chunking to the premise and aggregate predictions across chunks (Scirè et al., 2024; Yang et al., 2024). A sentence is counted as contradictory if the NLI model<sup>4</sup> predicts a contradiction label. The contradiction rate for a summary is defined as the fraction of its perturbed sentences labeled as contradictory, and dataset-level results are obtained by averaging these rates across summaries. The resulting contradiction rates for each dataset and perturbation type are reported in Table 3. Illustrative examples of original and perturbed summaries for selected perturbation types are shown in Figures 5, 6, & 7.

Across most perturbations, contradiction rates remain low, indicating that paraphrasing, simplification, vocabulary reduction, summarization, and source text insertion generally preserve factual consistency with respect to the original summaries. This supports our assumption that score changes observed in the main experiments primarily reflect metric sensitivity to surface and structural variation, rather than systematic factual errors introduced by the perturbations.

<sup>4</sup><https://huggingface.co/MoritzLaurer/DeBERTa-v3-large-mnli-fever-anli-ling-wanli>

We observe higher contradiction rates for the *Negated* perturbation across all datasets. This behavior is expected and does not necessarily indicate that these perturbations introduce factual errors. The *Negated* perturbation is explicitly designed to negate some statements in the summary, and therefore, a non-zero rate of contradictions indicates that the perturbation is being applied as intended. Importantly, the NLI-based validation used here is not designed to distinguish between *intended* negations that preserve overall factual meaning and *undesired* negations that fundamentally alter the factual content of the summary. Determining whether a negation invalidates the summary would require verifying each negated statement against the original source document, which in turn would necessitate fine-grained human evaluation over long inputs. Such an analysis is substantially more expensive and complex and lies beyond the scope of this work. As a result, higher contradiction rates for negated summaries should be interpreted as evidence that the perturbation successfully introduces negation, rather than as definitive proof of factual inconsistency. This limitation further highlights the need for more nuanced evaluation methods, including human verification, when assessing logical transformations in long-document summarization.

## D Per-Dataset Results

Tables 4, 5, and 6 report mean factuality scores for each metric under all seven perturbation types and for the original summaries across the three datasets used in this study. These tables provide the complete quantitative results corresponding to the aggregate trends shown in Figure 2. Consistent with our main analysis, MiniCheck and UniEval appear to be the most robust overall, maintaining relatively stable scores across most perturbations and datasets, with the exception of degraded performance on *Negated* summaries. In contrast, the remaining metrics are influenced by almost all types of perturbations, showing greater score variability, particularly in the legal domain.

BARTScore<sup>5</sup> performs poorly even on the original (unperturbed) summaries across all datasets. As a generation-based metric, its scoring depends heavily on the size and structure of the retrieved context, which may not be the ideal case in long-document setting, where evidence for a single sum-

<sup>5</sup>We use the implementation provided by Bishop et al. (2024).

**P1. Paraphrased**

system\_prompt: Provide the paraphrased version of the text.\n\nYou are strictly prohibited from omitting any information or altering its original meaning. Do not include explanations, reasoning, or commentary in your output.

user\_prompt: Text: <summary>

**P2. Less Diverse**

system\_prompt: Rewrite the following text using less diverse vocabulary.\n\nYou are strictly prohibited from omitting any information or altering its original meaning. Do not include explanations, reasoning, or commentary in your output.

user\_prompt: Text: <summary>

**P3. Negated**

system\_prompt: Rewrite the following text by introducing logically equivalent negations while preserving its original meaning.\n\nYou are strictly prohibited from omitting any information or altering its original meaning. Do not include explanations, reasoning, or commentary in your output.

user\_prompt: Text: <summary>

**P4. Simplified**

system\_prompt: Rewrite the following text by making complex sentences simpler.\n\nYou are strictly prohibited from omitting any information or altering its original meaning. Do not include explanations, reasoning, or commentary in your output.

user\_prompt: Text: <summary>

**P5. Summarized**

system\_prompt: Rewrite the text to make it more concise.\n\nYou are strictly prohibited from omitting any information or altering its original meaning. Do not include explanations, reasoning, or commentary in your output.

user\_prompt: Text: <summary>

**P6. Synonym Replacement**

system\_prompt: Revise the text using synonyms for some common words.\n\nYou are strictly prohibited from omitting any information or altering its original meaning. Do not include explanations, reasoning, or commentary in your output.

user\_prompt: Text: <summary>

**P7. Added Source Text**

system\_prompt: Insert a source sentence into the summary that does not relate to its main ideas.\n\nDo not include explanations, reasoning, or commentary in your output.

user\_prompt: Text: <summary> \n\n Source: <document>

Figure 4: Prompt templates used with GPT-4o to generate meaning-preserving perturbations of the original summaries.

Dataset	#Examples (Used)	Avg. Summary Sentences	Avg. Summary Tokens	Avg. Document Sentences	Avg. Document Tokens	Summary Type
SQuALITY	260	12.5	273	456.6	6,131	Human-written
LexAbSumm	351	4.2	169	385.9	10,840	Human-written
ScholarQABench	100	43.2	1,158	575.4	14,652	Human-written

Table 2: Dataset statistics for the three long-document summarization benchmarks used in this study.

1144 mary sentence may be scattered across distant sections of the source. The mismatch between the  
 1145 localized retrieved snippets and the broader document context can distort likelihood estimates, and  
 1146 this effect compounds when scores are aggregated over full summaries, leading to consistently lower  
 1147 values. We also observe that while a few individual sentences receive high BARTScore values, most  
 1148 have extremely low scores due to being more compressed and contextually demanding, which drives  
 1149 the overall average down.  
 1150  
 1151  
 1152  
 1153  
 1154

## 1155 E Domain-Specific Robustness

1156 To further analyze how text characteristics across domains influence factuality metric robustness, we  
 1157

1158 compute the mean absolute score change under meaning-preserving perturbations. While signed  
 1159 average score differences are often close to zero due to cancellation effects, absolute changes capture  
 1160 the magnitude of metric instability regardless of direction. For a given domain, metric, and perturbation,  
 1161 we compute:  
 1162  
 1163  
 1164

$$\Delta_{\text{abs}} = \frac{1}{N} \sum_{i=1}^N \left| M_{\text{pert}}^{(i)} - M_{\text{orig}}^{(i)} \right|, \quad (5)$$

1165 where  $M_{\text{pert}}^{(i)}$  and  $M_{\text{orig}}^{(i)}$  denote the factuality scores for the original and perturbed summaries  
 1166 of example  $i$ , respectively, and  $N$  is the number of summaries in the domain. This measure reflects the  
 1167  
 1168  
 1169

**Original summary:**

The story begins in thick jungle on Sekk, which we are told is a "second moon" which retains a "breathable atmosphere" around a lake surrounded by eleven jungled valleys. In this way, it is implied that Sekk is a second moon of Earth. In the jungle, we meet Noork and a young woman named Sarna. They begin traveling together through the jungle, but soon Sarna disappears and Noork is attacked. This is our first encounter with the Misty Ones, who blend in with the jungle foliage. Noork defeats the Misty Ones and continues toward the lake and island where they make their home. Noork briefly encounters his friend Ud near the marshy lowlands that lie between the jungled valleys on Sekk and the central Lake of Uzdon, but this area is not described. When Noork reaches the central island in the lake, we encounter a non-jungle landscape for the first time. Noork finds himself in a cultivated field, and sees the shape of a huge white skull about half a mile away. After speaking with an enslaved man and learning where Sarna is being held, Noork continues toward the skull. The skull is a dome of white stone, with black stone for eye-sockets and nose-holes. The interior contains a raised altar made of precious metals—gold, silver, and brass—and precious stones, as well as stone images of the two gods the Misty Ones worship. Below the altar is the caged area where the young women are held; Noork detects the entrance to this area by its foul odor. The room where the young women are kept is dimly lit by only two torches, very damp with pools of dirty water all around, and holds at least twenty young women. They have nothing to sit on but rotten grass mats. In contrast to the enslaved men who are out in the cultivated fields and open air, the young women are in a desperate situation indeed. They can only sit in their foul, rotting prison and wait to be sacrificed.

**Perturbed summary (Summarized):**

The story begins on Sekk, a second moon with a breathable atmosphere around a lake and eleven jungled valleys. Noork meets Sarna in the jungle, but she soon disappears, and he is attacked by the Misty Ones, whom he defeats. Noork encounters his friend Ud near marshy lowlands before reaching the central island, where he finds a cultivated field and a large white skull. After speaking with an enslaved man about Sarna's location, Noork approaches the skull, which is a dome of white stone with a raised altar inside, containing precious metals and images of the Misty Ones' gods. Below the altar, he discovers a dimly lit, damp room holding at least twenty young women, who are imprisoned and awaiting sacrifice.

Figure 5: Example of an original summary from SQUALITY and its Summarized perturbation.

**Original summary:**

The Court held that there had been a violation of Article 9 of the Convention, finding that the applicant did not enjoy the necessary procedural safeguards in having his request for alternative civilian service examined. The Court considered in particular that the Greek authorities had failed in their duty to ensure that the interviewing of conscientious objectors by the army's Special Board took place in conditions that guaranteed procedural efficiency and the equal representation required by domestic law. In this respect, it noted that: the applicant had been interviewed by a Board made up primarily of servicemen, two of the civilian members of the Board being absent but not replaced; the Minister of Defence's final decision, on the basis of a draft ministerial decision following the Board's proposal, did not afford the requisite safeguards of impartiality and independence; the scrutiny of the Supreme Administrative Court concerned only the lawfulness of the decision, not the merits, and was based on the assessments of the Special Board.

**Perturbed summary (Simplified):**

The Court found that Article 9 of the Convention was violated. The applicant did not have the necessary procedural safeguards for his request for alternative civilian service. The Court noted that the Greek authorities did not ensure that the army's Special Board interviewed conscientious objectors in a way that met procedural efficiency and equal representation as required by law. The applicant was interviewed by a Board mostly made up of servicemen, and two civilian members were absent without replacement. The Minister of Defence's final decision was based on a draft from the Board and did not ensure impartiality and independence. The Supreme Administrative Court only reviewed the lawfulness of the decision, not its merits, and relied on the Special Board's assessments.

Figure 6: Example of an original summary from LexAbSumm and its Simplified perturbation.

**Original summary:**

Solving large complex partial differential equations (PDEs), such as those that arise in computational fluid dynamics (CFD), is a computationally expensive process. This has motivated the use of deep learning approaches to approximate the PDE solutions. Convolutional neural networks (CNNs) have been shown to achieve a speedup of between two and four orders of magnitude over numerical flow solvers [40-42]. Data-driven neural solvers learn from observations of fluid dynamics, usually from large datasets generated by numerical simulation, without (in general) any knowledge of the governing equations. The type of DL model employed for data-driven solvers primarily depends on the nature of the physical problem, with CNNs being a natural choice for simulations where the geometry aligns with the coordinate system, while GNNs are preferred for simulations of more complex geometries and Lagrangian systems. More recently, physics-informed neural networks (PINNs) have been developed to solve forward and inverse problems where full or partial knowledge of the governing equations is known [23,24,29,30,103,104].

### 1. High-Dimensionality and Scale\*\*Problem:\*\* Fluid simulations often deal with high-dimensional spaces, particularly in three-dimensional simulations. Neural networks used for such tasks need to cater to vast input and output data scales.\*\*Impact:\*\* Handling high-dimensional input data while ensuring computational tractability remains a significant obstacle. The curse of dimensionality can lead to increased model complexity and training data requirements.

### 2. Generalization and Extrapolation\*\*Problem:\*\* Neural networks often struggle to generalize beyond the training data and may fail to extrapolate well to unseen conditions, which is problematic for fluid simulations that often involve varied and complex domains.\*\*Impact:\*\* Ensuring reliable performance across different fluid flows, geometries, and boundary conditions remains unsolved

### 3. Data Efficiency and Scarcity\*\*Problem:\*\* Generating the high-fidelity simulation data needed to train neural networks can be prohibitively expensive and time-consuming.\*\*Impact:\*\* NNs require large amounts of training data to generalize well. Data-efficient learning methods are critically needed to make the neural network approach feasible for fluid simulations.

### 5. Interpretability and Physics Consistency\*\*Problem:\*\* Neural networks often function as black boxes, providing little insight into how they arrive at their solutions or maintaining physical constraints.\*\*Impact:\*\* Understanding and ensuring that neural network predictions adhere to the underlying physical laws described by the PDEs is crucial for their effective application in fluid simulations.

**Perturbed summary (Negated):**

Solving large complex partial differential equations (PDEs), such as those that arise in computational fluid dynamics (CFD), is not a computationally inexpensive process. This has not discouraged the use of deep learning approaches to approximate the PDE solutions. Convolutional neural networks (CNNs) have not been shown to achieve a speedup of less than two and four orders of magnitude over numerical flow solvers [40-42]. Data-driven neural solvers do not learn from observations of fluid dynamics, usually from small datasets generated by numerical simulation, with (in general) some knowledge of the governing equations. The type of DL model employed for data-driven solvers does not primarily depend on the nature of the physical problem, with CNNs not being a natural choice for simulations where the geometry does not align with the coordinate system, while GNNs are not preferred for simulations of simpler geometries and Lagrangian systems. More recently, physics-informed neural networks (PINNs) have not been developed to solve forward and inverse problems where no knowledge of the governing equations is known [23,24,29,30,103,104].

### 1. High-Dimensionality and Scale\*\*Problem:\*\* Fluid simulations do not often deal with low-dimensional spaces, particularly in two-dimensional simulations. Neural networks used for such tasks do not need to cater to small input and output data scales.\*\*Impact:\*\* Handling low-dimensional input data while ensuring computational intractability does not remain a significant obstacle. The curse of dimensionality does not lead to decreased model complexity and training data requirements.

### 2. Generalization and Extrapolation\*\*Problem:\*\* Neural networks do not often succeed in generalizing beyond the training data and may succeed in extrapolating well to seen conditions, which is not problematic for fluid simulations that do not often involve varied and complex domains.\*\*Impact:\*\* Ensuring unreliable performance across different fluid flows, geometries, and boundary conditions remains solved.

### 3. Data Efficiency and Scarcity\*\*Problem:\*\* Generating the low-fidelity simulation data needed to train neural networks can be prohibitively inexpensive and time-saving.\*\*Impact:\*\* NNs do not require small amounts of training data to generalize poorly. Data-inefficient learning methods are not critically needed to make the neural network approach infeasible for fluid simulations.

### 5. Interpretability and Physics Consistency\*\*Problem:\*\* Neural networks do not often function as transparent boxes, providing much insight into how they arrive at their solutions or maintaining physical constraints.\*\*Impact:\*\* Understanding and ensuring that neural network predictions do not adhere to the underlying physical laws described by the PDEs is not crucial for their ineffective application in fluid simulations.

Figure 7: Example of an original summary from ScholarQABench and its Negated perturbation.

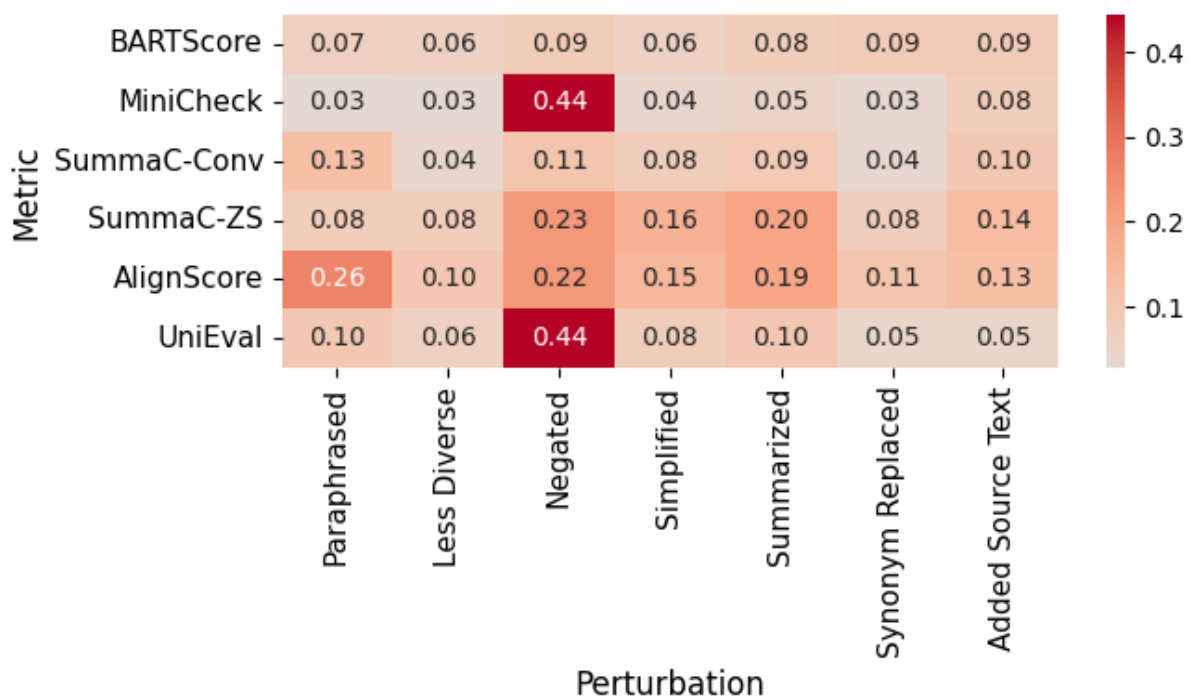


Figure 8: LexAbSumm: Metric score deltas under perturbations.

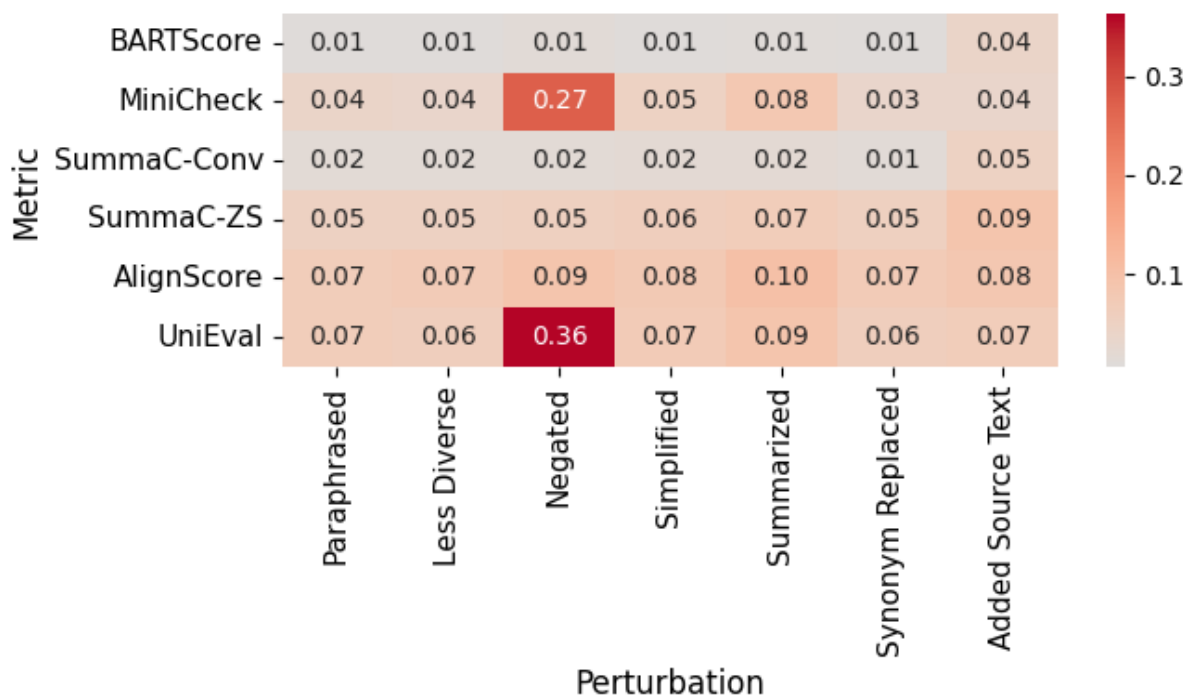


Figure 9: SQuALITY: Metric score deltas under perturbations.

Dataset	Paraphrased	Less Diverse	Negated	Simplified	Summarized	Synonym Replaced	Added Source Text
SQuALITY	0.023	0.033	0.681	0.029	0.018	0.034	0.052
LexAbSumm	0.004	0.001	0.560	0.002	0.007	0.013	0.030
ScholarQABench	0.010	0.017	0.542	0.013	0.006	0.028	0.019

Table 3: Average contradiction rate between perturbed and original summaries, computed using an NLI-based faithfulness check. Lower values indicate higher factual consistency.

Metric	Original	Synonym Replaced	Summarized	Simplified	Paraphrased	Negated	Less Diverse	Added Source Text
BARTScore	0.16	0.07	0.08	0.11	0.09	0.07	0.10	0.23
MiniCheck	0.84	0.83	0.84	0.85	0.85	0.40	0.85	0.78
SummaC-Conv	0.33	0.31	0.29	0.37	0.21	0.23	0.34	0.42
SummaC-ZS	0.38	0.39	0.32	0.47	0.39	0.22	0.38	0.48
AlignScore	0.52	0.48	0.42	0.56	0.38	0.38	0.51	0.58
UniEval	0.82	0.83	0.80	0.84	0.86	0.39	0.82	0.83

Table 4: Mean factuality scores for each metric and perturbation type on LexAbSumm.

Metric	Original	Synonym Replaced	Summarized	Simplified	Paraphrased	Negated	Less Diverse	Added Source Text
BARTScore	0.03	0.02	0.02	0.02	0.02	0.02	0.03	0.06
MiniCheck	0.32	0.32	0.32	0.35	0.32	0.19	0.35	0.39
SummaC-Conv	0.23	0.23	0.23	0.24	0.23	0.22	0.23	0.27
SummaC-ZS	0.18	0.19	0.19	0.22	0.18	0.13	0.19	0.23
AlignScore	0.21	0.22	0.19	0.22	0.20	0.20	0.22	0.26
UniEval	0.73	0.73	0.72	0.72	0.73	0.32	0.71	0.73

Table 5: Mean factuality scores for each metric and perturbation type on ScholarQABench.

Metric	Original	Synonym Replaced	Summarized	Simplified	Paraphrased	Negated	Less Diverse	Added Source Text
BARTScore	0.03	0.02	0.02	0.03	0.03	0.02	0.03	0.07
MiniCheck	0.56	0.55	0.53	0.55	0.56	0.30	0.55	0.57
SummaC-Conv	0.23	0.22	0.22	0.23	0.22	0.22	0.22	0.27
SummaC-ZS	0.13	0.14	0.10	0.14	0.12	0.11	0.12	0.19
AlignScore	0.18	0.20	0.11	0.16	0.16	0.20	0.15	0.22
UniEval	0.68	0.71	0.67	0.65	0.72	0.32	0.65	0.70

Table 6: Mean factuality scores for each metric and perturbation type on SQuALITY.

average magnitude of score variation induced by perturbations and serves as a robustness diagnostic.

In the legal domain (*LexAbSumm*), we observe the largest overall instability across both metrics and perturbations, as shown in Figure 8. Among the evaluated metrics, *AlignScore* exhibits the highest mean absolute score change, followed by *SummaC-ZS* and *UniEval*, indicating heightened sensitivity to surface-level changes in legally structured text. *MiniCheck* and *SummaC-Conv* show comparatively lower instability, while *BARTScore* exhibits the smallest absolute changes, consistent with its generally low baseline scores on this dataset. Across perturbations, *Negated* summaries produce by far the largest absolute score changes, followed by *Summarized* and *Paraphrased* variants. This pattern reflects the reliance of legal summaries on precise logical structure and domain-specific terminology, where even meaning-preserving changes can substantially alter cues used by factuality met-

rics.

For the narrative domain (*SQuALITY*), absolute score changes are smaller overall than in *LexAbSumm* but remain non-trivial, as shown in Figure 9. *UniEval*, *AlignScore*, and *MiniCheck* display the highest instability, while *SummaC-Conv* and *BARTScore* are comparatively stable. *Negated* summaries have the highest score change, and perturbations that increase abstraction, particularly *Summarized* and *Added Source Text*, induce the largest score changes. This suggests that narrative summaries, which often compress temporal and causal structure, pose challenges for metrics when abstraction increases, even if factual content is preserved.

In the scientific multi-document domain (*ScholarQABench*), we observe the lowest absolute score changes across all metrics and perturbations, as illustrated in Figure 10. Although *UniEval* and *MiniCheck* still exhibit measurable sensitivity, the

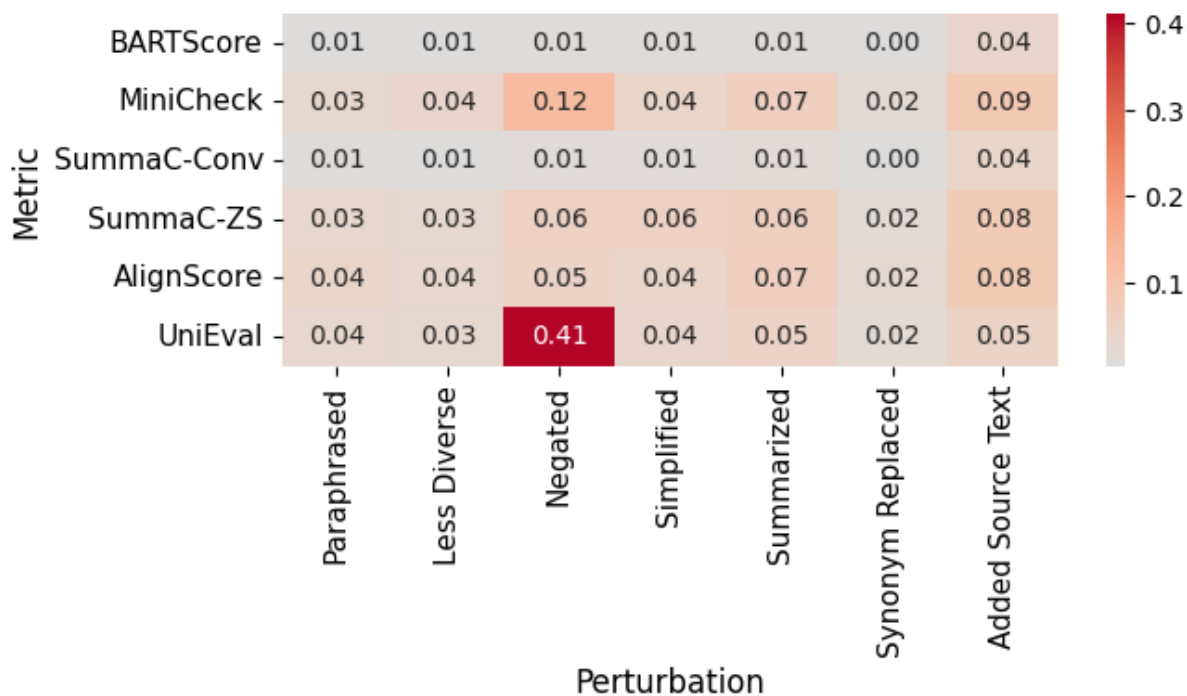


Figure 10: ScholarQABench: Metric score deltas under perturbations.

magnitude of instability is consistently lower than in the single-document domains. *Negated* and *Added Source Text* perturbations remain the most impactful, but their effects are attenuated. This relative stability likely arises from redundancy across multiple documents, where repeated evidence reduces the impact of localized reformulations on factuality assessment.

Overall, this analysis shows that factuality metric robustness varies substantially by domain, even under meaning-preserving perturbations. Legal text amplifies metric instability, narrative text exhibits moderate sensitivity to abstraction, and multi-document scientific text provides a stabilizing effect. These domain-specific patterns help explain the wide score distributions observed in Figure 2 and reinforce the need to evaluate factuality metrics under realistic long-document conditions.

## F Code & Data Availability Statement

We release all perturbed summary data, along with the recipes used to generate it and the scripts required to reproduce our results, under the MIT License. The complete codebase and dataset artifacts will be made publicly available upon publication.

## G Model Size & Budget

We describe all factuality metrics and the experimental setup in § 4. All experiments are conducted using a single NVIDIA H100 SXM5 80GB GPU.

## H Software Package Parameters

- NLTK (Bird, 2006): We use the punkt sentence tokenizer for sentence tokenization.
- OpenAI GPT-4o: We use top  $p$  sampling at 50% with a temperature of 0 for all prompts.