

LCFO: Long Context and Long Form Output Dataset and Benchmarking

Anonymous ACL submission

Abstract

This paper presents the Long Context and Form Output (LCFO) benchmark, a novel evaluation framework for assessing gradual summarization and summary expansion capabilities across diverse domains. LCFO consists of long input documents (5k words average length), each of which comes with three summaries of different lengths (20%, 10%, and 5% of the input text), as well as approximately 15 questions and answers (QA) related to the input content. Notably, LCFO also provides alignments between specific QA pairs and corresponding summaries in 7 domains.

The primary motivation behind providing summaries of different lengths is to establish a controllable framework for generating long texts from shorter inputs, i.e. summary expansion. To establish an evaluation metric framework for summarization and summary expansion, we provide human evaluation scores for human-generated outputs, as well as results from various state-of-the-art large language models (LLMs).

GPT-4o-mini achieves best human scores among automatic systems in both summarization and summary expansion tasks ($\approx +10\%$ and $+20\%$, respectively). It even surpasses human output quality in the case of short summaries ($\approx +7\%$). Overall automatic metrics achieve low correlations with human evaluation scores (≈ 0.4) but moderate correlation on specific evaluation aspects such as fluency and attribution (≈ 0.6).

1 Introduction

Robust long text generation capabilities are required to meet user demand for extensive content creation, including story writing and essay composition (Xie and Riedl, 2024), which is why recent models such as GPT-4 (OpenAI et al., 2024) are expanding the output lengths from 4k tokens in GPT-4o to 64k in the latest versions.

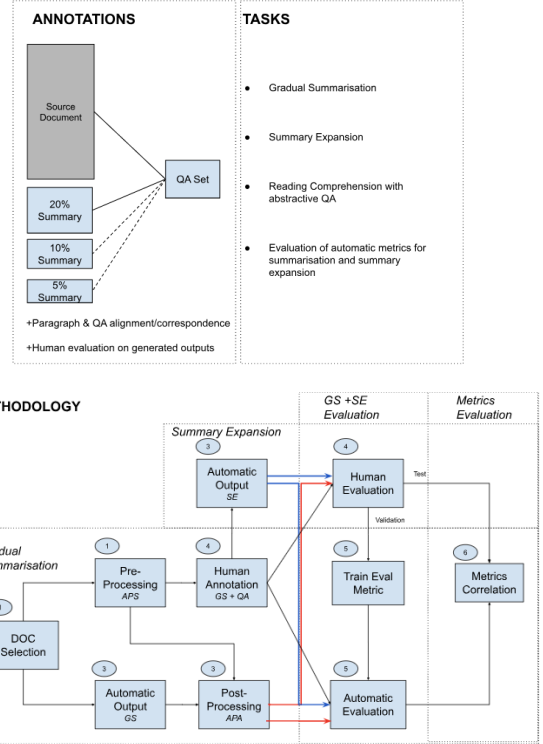


Figure 1: LCFO (top): annotations consist of 3 gradual summaries plus QA on each input document plus human evaluation annotations for gradual summarization and summary expansion; LCFO can be used for tasks such as gradual summarization, reading comprehension, summary expansion, and automatic metric evaluation. LCFO methodology (bottom) consists of 6 big steps.

However, evaluating the performance of Large Language Models (LLMs) on summarization and summary expansion tasks (see definitions in Section 2.1) is particularly challenging, especially when it comes to summarizing very long input documents and generating either long summaries or long summary expansions. Although there is a lot of work and interest in studying summarization evaluation (e.g. Zhang et al., 2024a), the evaluation of long text outputs is an emerging area (Que et al., 2024). Indeed, long-context input processing tasks

(such as summarization or comprehension question answering applied to long documents) and long-form output production both involve high cognitive loads for humans. This may be why evaluation work in these areas may not be as mature as in others.

The dataset presented here is an exclusively human annotated benchmark and challenge dataset involving natural language understanding and generation across multiple domains in the aforementioned tasks. Our data set is carefully manually crafted with human revision from the selection of the documents to the final annotation, without relying on LLMs at any point. We provide detailed linguistic guidelines and abstractive QA. Furthermore, we provide another set of linguistic guidelines to evaluate the tasks of summarization and summary expansion.

Together, we present the LCFO benchmark. Given a source of long structured documents, we generate multiple long outputs and associated QA pairs, and we evaluate human and model outputs with human annotations. The main contributions of this work are described below (see Figure 1-left for a schematic representation of the benchmark and its tasks):

- Dataset creation with structured inputs and alternative references; each input document is associated with 3 summaries of different lengths (20%, 10%, 5%). Gradual summarization is useful in that it provides both long and short summaries references. Moreover, the availability of summaries of various length is expected to improve the development of summary expansion approaches, which allows us to provide a more controllable summary expansion framework.
- A set of QA pairs for each input document aligned with each of the different length summaries. Our QA pairs are in free-form short-answer format (i.e. not multiple choice) and are of the abstractive type (i.e. they are not copies of parts of the source document). This QA can potentially be used to evaluate the model outputs, based on the appropriateness of the responses, similarly to previous proposals (Wang et al., 2022).
- Selection of automatic metrics. Most of these metrics can be used to evaluate at the paragraph level, and more widely can be used

to evaluate summary expansion or long-form generation in multiple tasks and languages.

- Evaluation of several LLMs on our dataset both automatically and manually; and evaluation of automatic metrics on summarization and summary expansion.

2 The LCFO Creation Framework: Principles and Methods

This section describes the detailed principles and methods that underlie LCFO and is shown in Figure 1 (right).

2.1 Definitions

Long context/form. We define long text as text that exceeds 5k words. For reference, the attention span for reading and taking notes has been considered 10 to 15 minutes since a 1978 seminal paper (Hartley and Davies, 1978). More recent research cited in a survey paper (Bradbury, 2016) shows that attention paid to lectures declines significantly after 20 minutes. The tasks we describe here are closer to note taking than lecture attendance; therefore, we should keep the 10–15 minute reference. If we base ourselves on the reported reading average time for first-language, secondary-education level readers (125 words per minute on average), we can consider that we will reach high cognitive load at around 1.5k words, and peak cognitive load at around 2.5k words. The quality of cognitive processing then starts to decrease significantly after 2.5 k words.

Structured/hierarchical. The input and output documents are partitioned into sections and, if necessary, nested sub-sections. Their structure is determined by task- and domain-specific guidelines. The size of the partitions is such that it makes them shorter (see examples of such structure in the data sets details).

Gradual summarization (GS). The input is a long document (defined here as 5k words or more). The summarization task can be described as the act of generating a much shorter corresponding document that includes the essential information contained in the long document and the same logical structure linking the various pieces of essential information. The summarization task presented here consists in taking long documents as inputs and generating three corresponding summaries of

length that represent 20%, 10%, or 5% of the input document.

Summary expansion (SE). The input is a short and concise document that has similar properties to those of a summary (that is, it is mainly a standalone document that abstracts from details). The summary expansion task can be described as the act of generating a much longer document that preserves the essential elements found in the corresponding short document as well as the logical structure that connects such elements. More specifically, the task presented here consists in taking summaries as inputs and generating 3 long documents of different lengths. Each of the 3 lengths is set such that an input summary represents either 5%, 10%, or 20% of its respective expanded documents. As this is a more freely generative task, an additional requirement to be taken into consideration is that of coherence (for example, the detailed information included in one generated sentence should not contradict that included in another sentence).

2.2 Data selection and preprocessing

Selection. We select input documents that cover different domains, which meet the desired average length and the structure requirement.

- We cover 7 domains including politics, news, Wikipedia, scientific, literature, conversational, and legal documents, with the format of documents and conversations. We source from 10 datasets: LexGLUE (Chalkidis et al., 2022), BookSum (Kryscinski et al., 2022), SQuality (Wang et al., 2022), FacetSum (Meng et al., 2021) JRC-Acquis (Steinberger et al., 2006), MultiUN, Wikipedia, GovReport (Huang et al., 2021), Summscreen (Chen et al., 2022), and Seahorse (Clark et al., 2023). The correspondence between the source data sets and the domains is shown in Table 1.
- The source documents are selected to be on average 5k words / document. We prefer documents with a hierarchical structure and containing relatively few numbers, which are better suited for summarization and summary expansion tasks.
- We prioritize recent documents when the domains allow (e.g., Wikipedia, where articles have a significant amount of new information

since 2024). We preprocess documents in structured domains to provide a flattened structure while keeping hierarchical markers that are readable to annotators and models. It also ensures a consistent format across datasets.

- We also filter out documents that contain toxicity using the ETOX package (Costa-jussà et al., 2023) and add manual verification to ensure the high integrity of the selected documents.

Preprocessing. To reduce the cognitive load for human annotation, we split paragraphs automatically (APS). Details on this paragraph splitting differ from corpus to corpus and are reported in the Appendix A.

2.3 Human summary and QA-pair generation

To obtain human-written summaries of long-form texts, detailed guidelines are developed, which are shown in Appendix D. All summary writers must be native English speakers and have writing or editing experience.

These writers receive 252 long form documents (each around 5k words), and they are asked to read each document in its entirety and write three summaries for each document: the first summary representing around 20% of the length of the source document; the second and the third summaries representing further summarization — around 10% and 5% of the source document, respectively. When writing these summaries, the writers are tasked with compressing and retaining all the core ideas of the source. Giving a definition of a “core” or “main” idea of a text presented one of the challenges of our work with the writers. Each summary is supposed to be a cohesive standalone text that could be read and understood on its own.

The fact that the source documents represent different domains poses another challenge for the writers: they need to possess some knowledge and expertise in each of these. The documents are split into sections and paragraphs, and the writers are asked to keep the flow of the section/paragraph structure of the source text, trying to summarize it from top to bottom. However, we emphasize that the sentence by sentence summarization is exactly what we do not need, and the summaries need to be abstractive rather than extractive, which means copying the source text is strongly discouraged.

In addition, the writers are asked to provide a set of questions and answers for each long document.

They need to compose at least 13–15 questions per 5,000 words. The answers are supposed to cover the points reflected in the summaries. We instruct the writers to produce open-ended, complex questions, which provide a good baseline for testing reasoning. For tracking and alignment purposes, each paragraph in each source document is given a number. The writers are asked to specify in which paragraph each answer can be found. They are also asked to indicate which of the summaries provides the answer to each question.

Besides the general guidance, we discovered that working with conversational content needed additional clarifications, so we prepared an additional document for working specifically with long-form text that contains conversations (such as plays or screenplays).

2.4 Automatic output and postprocessing

We want to understand how current state-of-the-art models perform on our new benchmark, both on the capability of comprehension a very long context and on the generation of long outputs. We conduct the automatic abstractive summarization for the former and summary expansion for the latter. We give details on the tasks below.

Gradual Summarization. We prompt the models with the human guidelines with a slight adaptation to be LLM friendly. We input the entire document without paragraph splitting. To give a fair evaluation, we prompt all LLMs in the zero-shot setting. To control the length of the LLM output, we have added additional instructions with the upper and lower bounds of the permissible words. For example, to ask the model to generate a summary of the $R\%$ length of the source text, the prompt contains "Make sure the summary has $\{y\}$ words or less.....Please write at least $\{x\}$ words""", where x and y are determined per document with respect to the length and ratio R . In practice, we see that enforcing the length of the document right before and after the content block in the prompt gives consistent results. We give details of our summary prompts in the Appendix H.

Summary expansion. We customize prompts for each domain, plus the model-specific prompt templates. Similarly to the summarization task, the prompt contains instructions on the desired range of the generated text length. In addition, each prompt has instructions to guide the model in generating

content of a certain quality (consistency, coherence, and keeping the main ideas in the summary). We prompt the model with specific formats for different domains as reported in the appendix H.

Automatic Paragraph Alignment (APA). We add this step because we want to perform a human evaluation of the outputs that we are creating (as detailed in the next section 2.5). The task of comparing long inputs and outputs creates a high cognitive load on human evaluators. To reduce it, we provide an approximate alignment between the input paragraphs and the segments of the output, taking advantage of the assumption that a summary usually follows the structure of the source document. First, we use dynamic programming to find a monotonic alignment path between input and output sentences that would maximize the sum of cosine similarities of the SONAR embeddings (Duquenne et al., 2023) of the two sentences. An output sentence could be aligned with multiple consecutive input sentences, potentially from different paragraphs, but we assign it to a single input paragraph with which it is aligned the most frequently. In this way, each input paragraph gets aligned to a contiguous output segment (potentially empty) in a monotonic way. This alignment helps the annotators navigate the input and output documents in a joint way.

2.5 Human Evaluation

To perform human evaluation on previously generated output, we design human evaluation guidelines inspired by previous works (Clark et al., 2023; Krishna et al., 2023; Que et al., 2024) and fully reported in Appendices F and G.

Human evaluation on gradual summaries. Before starting the evaluation, annotators are allowed to reject a task if the output text is gibberish or obviously of low quality.

The generated summaries are evaluated in two tasks. In Task 1, the annotators first read the source document and the three summaries and then rate the generated text in four aspects, including attribution, coverage of the main ideas, conciseness and readability (similar to the 'checklist' in HelloBench, (Que et al., 2024)). The annotators rate the summary on a 0-4 Likert scale and finally give an overall rating on a 0-10 Likert scale. Each summary receives its own separate set of scores.

In Task 2, the annotators validate the QA sets that were previously created by human writers. For

each question in the QA sets (13–15 questions and answers), the annotators are required to determine whether the content of the summary contains enough information to answer the question (i.e., the answer is directly stated, heavily implied or logically entailed in the summary). The annotators give a YES or NO to each question and answer. For each summary, the annotators validate the whole set of QA once.

The whole evaluation is referenceless, which means that the human written summaries are not shown to annotators, and that they only see a single set of summaries from one anonymous model output each time.

Both tasks 1 and 2 involve human judgment, and to reduce the bias, 3 sets of rating from random annotators are required for each generated output. The same guidance should be used for all different domains. Detailed evaluation guidelines are included in the Appendix F.

Human evaluation on summary expansion. We use the same format as the previous summarization evaluation tasks and integrate some of the questions from Story Plot Generator (Zhu et al., 2023) and HelloBench (Que et al., 2024). For task 1, the annotators read the source summary and the generated long-form output, and rate the generated text on six aspects, including the coverage of main core ideas, cohesion, richness in details, creativity, non-repetitiveness, and interest, then give an overall rating at the end. In task 2 they validate the QA set with the generated long-form text. Each output is evaluated separately without reference and three sets of random annotation ratings are required. Detailed evaluation guidelines are included in the Appendix G.

Evaluation statistics. Summaries and summary expansions are each evaluated separately. For the evaluation of generated summaries, 252 documents from all domains are used as the source to generate the summaries (with 2 documents being excluded during the process). The summaries are generated using three different models (as reported in Section 3 and chosen to represent close and open models of different sizes): GPT-4o-mini-64k (OpenAI et al., 2024), Llama-3.1-70B, and Llama-3.1-8B (et al., 2024). This results in 756 outputs and, along with 252 sets of human-written summaries, creates a dataset of 1,000 document-summary pairs for evaluation. A vendor sources 287 annotators, who are required (1) to be native speakers of English and (2)

to hold a language-related degree. These annotators are selected from a pool that is different from that of the summary writers, ensuring that they have no prior knowledge of the source documents or the written summaries. Tasks are randomly assigned to annotators until every set of generated output receives three complete annotations. A limit of 10 evaluations for each model is set per annotator to mitigate biases in the results.

For the evaluation of generated summary expansions, only a subset of data (only covering domains that do not include factual information) is selected for evaluation, including SummScreen, BookSum, SQuality and FacetSum (102 source documents in total). The expansions are generated with the same models as previously (GPT-4o-mini, Llama-3.1-70B, and Llama-3.1-8B), resulting in 306 long-form outputs. Ten experienced data analysts are selected to conduct the evaluation. Similarly to the evaluation of summaries, the tasks are assigned randomly until every long form output receives 3 complete annotations.

2.6 Automatic evaluation

The summarization outputs are typically evaluated by computing the Rouge scores (Lin, 2004) with respect to a reference. However, this approach is not sufficient for at least three reasons (Schluter, 2017): It depends too much on the reference, it offers only a comparison at the surface level, and it does not explain why a summary is good or bad. Thus, we compute several other reference-free metrics. Each targets a specific property:

1. **Repetitiveness:** how much the summary repeats the same phrases. We report the number of all the word n-grams (with n from 1 to 3) in the summary, divided by the number of such unique n-grams (REP-3) (Welleck et al., 2019).
2. **Fluency:** how grammatical the text is. We report the average probability of a summary sentence being grammatical (or linguistically acceptable in the Chomskyan sense of the term) computed with a CoLA classifier (Krishna et al., 2020).
3. **Coherence:** how similar are the sentences in the generated texts to each other. COH-2 averages similarity of the neighboring-over-one sentences (in the embedding space) Parola et al. (2023).

4. **Attribution:** how much of the summary is directly attributable to the source (something like “precision” of the ideas in the summary). The average score of SEAHORSE Q4 (SH-4) model evaluates attribution (Clark et al., 2023).
5. **Coverage of the source:** how much of the source is reflected in the summary. The average SEAHORSE Q5 (SH-5) score reports this aspect (Clark et al., 2023).
6. **Overall** in order to evaluate the overall quality of the text we use the following metrics:
 - the aggregated score (AVG) from the above metrics for summarization, namely, is the average score of $-RE-3, CoLA, COH-2, SH-4, SH-5$. For summary expansion, the aggregated score (AVG) is the average score of $REP-3, CoLA, COH-2$. Note that $REP-3$ is negated to make the score monotonic. Also, for the summary expansion, $REP-3$ increases the value over the length of output, so the factor 0.2 is empirically set to normalize the value on the 20% summary expansion task.
 - we use HelloEval (HE) score (Que et al., 2024), which is an LLM-as-judge model with different checklist trained in human evaluation.

For some model-based scores (SEAHORSE), we had to feed the whole source text to a transformer model, which was neither feasible computationally with long context inputs nor made sense given the relatively short-form training data of those models. To bypass this problem, we segment sources and summaries into aligned fragments (using a modification of the alignment algorithm in Section 2.4) with at most 50 sentences on the source side and compute model-based metrics for the segment pairs. Table 5 in Appendix C summarizes the list of metrics.

2.7 Data Statistics

LCFO covers 7 domains sourced from 10 datasets with an average document length of 5k words. Table 1 contains the distribution of the LCFO dataset in subsets and domains, as well as the average word length of the documents. More details are reported in the Appendix B.

3 Experiments

Settings. We experimented with closed and open LLMs. We chose GPT-4o-mini-64k¹ for the closed model and Llama-3.1-70B (Dubey et al., 2024) for the open-source one. For summarization, we ran the model with all length ratios (5%, 10%, 20%), while for summary expansion, we only expanded the summaries 20% to the full document. We also performed a postprocessing step to filter the templated response such as “**Summary**”, “Here is the summary:”, etc.

Summarization results. Table 2 shows the general results of the selected models at different levels of gradual summarization. Results broken down by domains are reported in the Appendix I. Note that LLMs tend to perform similarly regardless of the length of the output in terms of human scores. This is not the case for humans that show to lag behind when performing short summaries.

The best results are consistently achieved with GPT-4o-mini and are consistent with previous research findings (Que et al., 2024). This model even surpasses human-level quality in short summaries. This may be explained by humans tending to perform worse when summarizing short documents and better when summarizing long ones.

Summary expansion results. Table 3 shows the overall results on the summary expansion task only on a factor of 5. The performance of models is not coherent across metrics that look only at the output (i.e. %WC, REP-3, CoLA, COH-2, and AVG). In terms of HE and coherently with the human evaluation results, GPT-4o-mini is the best performing model.

We do not report results on expanding by larger factors (10 or 20) since models are performing poorly (see some examples in the Appendix I). The main limitation in this case is that the models do not reach the required output length.

When comparing across tasks (i.e. summarizing to 20% or doing summary expansion from 20% by a factor of 5), the results show better performance in the former. It is expected that summary expansion is a more challenging task across domains and all models. Current models struggle with this task. If we compare HE, the deltas in the same model vary from 6% for GPT-4o-mini to $\approx 30\%$ for Llama-3.1-70B. When comparing output-based

¹<https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>

DATASET	# DOCS	DOMAIN	AVG LEN DOC
LexGLUE	25	Legal: supreme court opinions	4953
BookSum	27	Literary: books, novel, act	4114
SQuality	25	Literature: stories	4856
FacetSum	25	Scientific: journal articles on various domains	4904
JRC-Acquis	25	Legal: legislative text of the European Union	4825
MultiUN	25	Political: UN docs	4539
Wikipedia	25	Wikipedia: 22 docs on biomedicine	5266
GovReport	25	Political: Congressional Research Service and US Government Accountability Office	5078
Summscreen	25	Conversational: TV series transcript	5030
Seahorse	25	News: English BBC news	4576
Total	252	average word count	4814

Table 1: LCFO Summary: Domains and Statistics

Output	R-L(↑)	REP-3(↓)	CoLA↑	COH-2↑	SH-4↑	SH-5↑	AVG↑	HE↑	Hum↑
LCFO.5%									
Human	n/a	0.308	0.941	0.809	0.644	0.387	0.494	52.195	6.61
GPT-4o-mini	0.331	0.328	0.968	0.719	0.635	0.487	0.496	76.917	7.25
Llama-3.1-70B	0.384	0.383	0.965	0.861	0.622	0.377	0.488	72.468	6.27
Llama-3.1-8B	0.377	0.411	0.969	0.865	0.618	0.372	0.482	63.894	6.32
LCFO.10%									
Human	n/a	0.395	0.945	0.816	0.661	0.416	0.489	64.688	7.44
GPT-4o-mini	0.385	0.404	0.964	0.695	0.621	0.471	0.469	77.863	7.50
Llama-3.1-70B	0.434	0.515	0.944	0.860	0.614	0.369	0.454	72.497	6.42
Llama-3.1-8B	0.429	0.534	0.963	0.858	0.612	0.366	0.453	59.385	6.63
LCFO.20%									
Human	n/a	0.244	0.938	0.805	0.615	0.357	0.494	69.745	7.78
GPT-4o-mini	0.445	0.497	0.961	0.673	0.616	0.464	0.443	76.706	7.52
Llama-3.1-70B	0.467	0.631	0.928	0.860	0.596	0.357	0.422	71.603	6.32
Llama-3.1-8B	0.469	0.647	0.956	0.861	0.594	0.370	0.427	51.015	6.60

Table 2: Performance on the summarization task

Output	%WC	REP-3(↓)	CoLA↑	COH-2↑	AVG↑	HE↑	Hum↑
GPT-4o-mini	1.931	0.707	0.913	0.609	0.460	70.896	6.431
Llama-3.1-70B	1.058	0.680	0.877	0.750	0.497	39.199	4.469
Llama-3.1-8B	1.187	0.809	0.903	0.779	0.507	38.416	4.801

Table 3: Performance on the summary expansion task by a factor of 5, giving the 20% summary input.

metrics, there are discrepancies in conclusions (i.e., Llama-3.1-8B better than 70B model). However, HE is still worse for the 8B model. This may indicate that selected output-based quality metrics are less reliable than the HE score (see the analysis below for metrics evaluation).

Metrics evaluation. In our study, we consider human evaluation, conducted according to the guidelines outlined in Appendices F and G, as the definitive measure of the overall quality score, as well as the scores for individual quality aspects such as coverage and attribution. To mitigate potential

biases among the annotators, we calculate the average of three annotations for each task.

Table 4 presents the Spearman correlation coefficients for various aspects and overall scores, comparing automatic metrics and human evaluations for summarization and summary expansion, respectively. A higher Spearman correlation coefficient signifies a stronger correlation between the automatic metrics and human annotation. The metric that shows the highest correlation with human annotation corresponds to SH-4, which measures attribution. When comparing metrics that measure overall performance, we observe that R-L is not

	R-L	CoLA	SH-4	SH-5	AVG	HE
S	0.196 (0.065)	0.595 (6.337e-10)	0.616 (1.005e-10)	0.445 (1.105e-5)	0.159 (0.135)	0.428 (2.591e-05)
SE	n/a	n/a	n/a	n/a	0.285 (3.646e-05)	0.405 (1.957e-09)

Table 4: Spearman correlation coefficients (and p-value) for various aspects and overall scores between automatic metrics and human evaluation for summarization (S) and summary expansion (SE). For the former, we show correlations between CoLA and Human evaluation Q2d; SH-4 and Human evaluation Q2a; SH-5 and Human evaluation Q2b and R-L/AVG/HE and Human evaluation Q3 (appendix F). For the latter, we show correlations between AVG/HE and Human evaluation Q3 (appendix G).

very good at correlation, but it may also be due to the fact that our task is not the best suited to use human references. HE is the one with the highest correlation. AVG low-correlation (both in S and SE) may be explained by the fact that individual averaged metrics are not very good or they cover more specific aspects which may not end capturing the overall performance. This low correlation for R-L and AVG can explain the discrepancy observed in the model ranking (specially between Llama-3.1-70B and 8B in Tables 2 and 3. Correlations are low in all cases, which shows the difficulty of the evaluation. Beyond the challenge of automatizing it, we should add the fact that humans struggle in generating short summaries, which may imply that humans also struggle in evaluating them.

4 Related Work

Related work on long context and long form output comes in many flavors. In this section, we cover a summary of related works on long context and long-form output datasets.

Long-context datasets. Infinite length datasets such as NIAH, RULER (Hsieh et al., 2024) work with distracting information. Finite-length nondistractive-based datasets include: Longbench (Bai et al., 2024) and Marathon (Zhang et al., 2024b) that includes tasks with 5–25k context and, more recently, (Kwan et al., 2024) build a dataset up to 8k tokens context length to evaluate LLMs’ long-context understanding across five key abilities: understanding of single or multiple relevant spans in long contexts based on explicit or semantic hints, and global context understanding. Loong is a multi-document QA dataset up to 200k context to assess RAG abilities. HelloBench (Que et al., 2024) includes summarization of a selection of long-input documents (3k to 6k word length).

Long-form output datasets. There is a lack of reference-based datasets on long form output. How-

ever, there are datasets that study prompting of different long-form generation; e.g., StoryGen (Zhu et al., 2023) includes prompts to generate stories and HelloBench is one of the most diverse text long form generation including stories, screenplays, key-word writing.

Our contribution on datasets involve the manual collection of 3-length summaries from long input documents. This collection also includes abstractive QA (non-multiple choice) to test comprehension. Our contribution on metrics involves new human evaluation protocols on summarization and summary expansion, as well as annotations to train/evaluate supervised metrics on long-form outputs.

5 Conclusions

LCFO provides gradual summaries references from 5k input documents with QA pairs for each of the documents and summaries. Additionally, we provide human evaluation of human and model-generated summaries and model-generated summary expansions. Overall, LCFO enables the evaluation of several tasks and metrics in the setting of long-context input documents and long-form output. While the main contribution of this paper is to present the freely available LCFO dataset², we also evaluate model and human outputs, showing that LLMs are capable of surpassing human results when producing short summaries. Current evaluation results question the usefulness of manually generating human references for short summarization of long documents. To confirm this, as further work, we plan to exploit the capabilities of LCFO by using QA as part of automatic evaluation (i.e. scoring how many questions are correctly answered in model-generated summaries).

²Available at BLIND

Limitations and Ethical considerations

Data contamination. Source documents may exist in the training data of the models, therefore, generation may be at risk. To mitigate this, we prioritize recent documents, since this is not enough, we annotate the correspondence of sections in summarization versions, so that we can generate only portions of the document. Therefore, if the model uses internal knowledge, we can quantify by spotting details from other sections.

Experiments. The experimental options that LCFO offers are much larger than the ones we explore in this paper. Also, the dataset can be easily expanded to have more summary references by matching with existing summaries in some of the domains. QA pairs have not been used in the paper but this is designed (but not limited) to serve for doing reading comprehension and/or for creating an evaluation metric.

Metrics. Summarization is a generative task with very diverse aspects of quality, and no single automatic evaluation metric captures them all adequately. To compensate for this, we report multiple evaluation metrics, but still, some of them are not well established; for example, there is no single metric of longform text coherence that the summarization community agrees upon. By providing the results of the human evaluation, we hope to help the community develop and validate better automatic evaluation metrics in the future.

Computing In terms of computing, evaluating LLMs on LCFO benchmark require larger memory due to both its big context size and the long-form output (should the models be capable to it). In case of Llama, we used 1 NVIDIA GPU A100 80 GB for the 8B model, and 8 GPUs for the 70B model. The resource was shared with the loading of scoring models (SH, CoLA) as well. Each evaluation run over 10 domains takes 90 minutes, including the computation of all the scores except HelloEval (where the computing time depends on the external availability of GPT-4 endpoint deployment).

Annotations Annotators were paid a fair rate. Each of the annotators signed a consent form agreeing on the dataset and its usage that they were participating in.

References

- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. [LongBench: A bilingual, multi-task benchmark for long context understanding](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3119–3137, Bangkok, Thailand. Association for Computational Linguistics.
- N. Bradbury. 2016. Attention span during lectures: 8 seconds, 10 minutes, or more? *Advances in Physiology Education*, 40(4):509–513.
- Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Katz, and Nikolaos Aletras. 2022. [LexGLUE: A benchmark dataset for legal language understanding in English](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4310–4330, Dublin, Ireland. Association for Computational Linguistics.
- Mingda Chen, Zewei Chu, Sam Wiseman, and Kevin Gimpel. 2022. [SummScreen: A dataset for abstractive screenplay summarization](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8602–8615, Dublin, Ireland. Association for Computational Linguistics.
- Elizabeth Clark, Shruti Rijhwani, Sebastian Gehrmann, Joshua Maynez, Roei Aharoni, Vitaly Nikolaev, Thibault Sellam, Aditya Siddhant, Dipanjan Das, and Ankur Parikh. 2023. [SEAHORSE: A multilingual, multifaceted dataset for summarization evaluation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9397–9413, Singapore. Association for Computational Linguistics.
- Marta Costa-jussà, Eric Smith, Christophe Ropers, Daniel Licht, Jean Maillard, Javier Ferrando, and Carlos Escolano. 2023. [Toxicity in multilingual machine translation at scale](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9570–9586, Singapore. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Paul-Ambroise Duquenne, Holger Schwenk, and Benoit Sagot. 2023. [SONAR: sentence-level multimodal and language-agnostic representations](#). *arXiv preprint*.
- Abhimanyu Dubey et al. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.

- James Hartley and Ivor K. Davies. 1978. Note-taking: A critical review. *Programmed learning and educational technology*, 15(3):207–224.
- Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekish, Fei Jia, Yang Zhang, and Boris Ginsburg. 2024. [Ruler: What’s the real context size of your long-context language models?](#) *Preprint*, arXiv:2404.06654.
- Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. 2021. [Efficient attentions for long document summarization](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1419–1436, Online. Association for Computational Linguistics.
- Kalpesh Krishna, Erin Bransom, Bailey Kuehl, Mohit Iyyer, Pradeep Dasigi, Arman Cohan, and Kyle Lo. 2023. [LongEval: Guidelines for human evaluation of faithfulness in long-form summarization](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1650–1669, Dubrovnik, Croatia. Association for Computational Linguistics.
- Kalpesh Krishna, John Wieting, and Mohit Iyyer. 2020. Reformulating unsupervised style transfer as paraphrase generation. In *Empirical Methods in Natural Language Processing*.
- Wojciech Kryscinski, Nazneen Rajani, Divyansh Agarwal, Caiming Xiong, and Dragomir Radev. 2022. [BOOKSUM: A collection of datasets for long-form narrative summarization](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6536–6558, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Wai-Chung Kwan, Xingshan Zeng, Yufei Wang, Yusen Sun, Liangyou Li, Yuxin Jiang, Lifeng Shang, Qun Liu, and Kam-Fai Wong. 2024. [M4LE: A multi-ability multi-range multi-task multi-domain long-context evaluation benchmark for large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15568–15592, Bangkok, Thailand. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Rui Meng, Khushboo Thaker, Lei Zhang, Yue Dong, Xingdi Yuan, Tong Wang, and Daqing He. 2021. [Bringing structure into summaries: a faceted summarization dataset for long scientific documents](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 1080–1089, Online. Association for Computational Linguistics.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rameev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-

866	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Di-	925
867	Girish Sastry, Heather Schmidt, David Schnurr, John	nan, Kyunghyun Cho, and Jason Weston. 2019. Neu-	926
868	Schulman, Daniel Selsam, Kyla Sheppard, Toki	ral text generation with unlikelihood training. <i>arXiv</i>	927
869	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	<i>preprint arXiv:1908.04319</i> .	928
870	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,		
871	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	Kaige Xie and Mark Riedl. 2024. Creating suspenseful	929
872	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	stories: Iterative planning with large language mod-	930
873	lipe Petroski Such, Natalie Summers, Ilya Sutskever,	els . In <i>Proceedings of the 18th Conference of the</i>	931
874	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	<i>European Chapter of the Association for Computa-</i>	932
875	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	<i>tional Linguistics (Volume 1: Long Papers)</i> , pages	933
876	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	2391–2407, St. Julian’s, Malta. Association for Com-	934
877	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	putational Linguistics.	935
878	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,		
879	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	Haopeng Zhang, Philip S. Yu, and Jiawei Zhang. 2024a.	936
880	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	A systematic survey of text summarization: From sta-	937
881	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	tistical methods to large language models . <i>Preprint</i> ,	938
882	Clemens Winter, Samuel Wolrich, Hannah Wong,	<i>arXiv:2406.11289</i> .	939
883	Lauren Workman, Sherwin Wu, Jeff Wu, Michael		
884	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	Lei Zhang, Yunshui Li, Ziqiang Liu, Jiayi Yang, Junhao	940
885	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	Liu, Longze Chen, Run Luo, and Min Yang. 2024b.	941
886	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	Marathon: A race through the realm of long con-	942
887	Zheng, Juntang Zhuang, William Zhuk, and Bar-	text with large language models . In <i>Proceedings</i>	943
888	ret Zoph. 2024. Gpt-4 technical report . <i>Preprint</i> ,	<i>of the 62nd Annual Meeting of the Association for</i>	944
889	<i>arXiv:2303.08774</i> .	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	945
		pages 5201–5217, Bangkok, Thailand. Association	946
		for Computational Linguistics.	947
890	Alberto Parola, Jessica Mary Lin, Arndis Simonsen,		
891	Vibeke Bliksted, Yuan Zhou, Huiling Wang, Lana	Hanlin Zhu, Andrew Cohen, Danqing Wang, Kevin	948
892	Inoue, Katja Koelkebeck, and Riccardo Fusaroli.	Yang, Xiaomeng Yang, Jiantao Jiao, and Yuan-	949
893	2023. Speech disturbances in schizophrenia: As-	dong Tian. 2023. End-to-end story plot generator .	950
894	sessing cross-linguistic generalizability of nlp au-	<i>Preprint</i> , <i>arXiv:2310.08796</i> .	951
895	tomated measures of coherence. <i>Schizophrenia Re-</i>		
896	<i>search</i> , 259:59–70.		
897	Haoran Que, Feiyu Duan, Liqun He, Yutao Mou,	A Automatic Paragraph Splitting Details	952
898	Wangchunshu Zhou, Jiaheng Liu, Wenge Rong,		
899	Zekun Moore Wang, Jian Yang, Ge Zhang, Junran	As part of the guidelines for summarization, anno-	953
900	Peng, Zhaoxiang Zhang, Songyang Zhang, and Kai	tators are instructed to read long documents from	954
901	Chen. 2024. Hellobench: Evaluating long text gener-	different domains and mentally distill key points	955
902	ation capabilities of large language models . <i>Preprint</i> ,	into paragraphs and then into a cohesive document	956
903	<i>arXiv:2409.16191</i> .	summary. This process requires the source text to	957
		be logically segmented into well-structured para-	958
		graphs that facilitate comprehension and synthesis.	959
904	Natalie Schluter. 2017. The limits of automatic sum-	During pilot studies, it became evident that the	960
905	marisation according to rouge. In <i>Proceedings of the</i>	quality of the initial paragraph segmentation sig-	961
906	<i>15th Conference of the European Chapter of the Asso-</i>	nificantly impacted annotation outcomes. Poorly	962
907	<i>ciation for Computational Linguistics</i> , pages 41–45.	segmented paragraphs increased cognitive load and	963
908	Association for Computational Linguistics.	risked misinterpretation, while cohesive and log-	964
909	Ralf Steinberger, Bruno Pouliquen, Anna Widiger,	ically structured paragraphs improved annotation	965
910	Camelia Ignat, Tomaž Erjavec, Dan Tufiş, and Dániel	consistency and efficiency.	966
911	Varga. 2006. The JRC-Acquis: A multilingual		
912	aligned parallel corpus with 20+ languages . In	Given the variability in formats and structures	967
913	<i>Proceedings of the Fifth International Conference</i>	across datasets, a uniform approach to paragraph	968
914	<i>on Language Resources and Evaluation (LREC’06)</i> ,	splitting was not feasible. Some datasets provided	969
915	Genoa, Italy. European Language Resources Associ-	explicit structural markers (e.g., new lines, section	970
916	ation (ELRA).	headers), while others required more algorithmic	971
917	Alex Wang, Richard Yuanzhe Pang, Angelica Chen, Ja-	intervention, such as employing the Segment Any-	972
918	son Phang, and Samuel R. Bowman. 2022. SQuAL-	thing Text (SaT-I3) model. Furthermore, the SaT-	973
919	ITY: Building a long-document summarization	I3 model’s performance varied across text types,	974
920	dataset the hard way . In <i>Proceedings of the 2022 Con-</i>	necessitating dataset-specific thresholds and post-	975
921	<i>ference on Empirical Methods in Natural Language</i>	processing techniques to optimize paragraph seg-	976
922	<i>Processing</i> , pages 1139–1156, Abu Dhabi, United	mentation.	977
923	Arab Emirates. Association for Computational Lin-		
924	guistics.		

978	This section outlines the tailored methodologies	1026
979	applied to each dataset in the LCFO corpus, high-	1027
980	lighting how their unique characteristics were ad-	1028
981	dressed to produce high-quality, preprocessed docu-	1029
982	ments for annotation.	1030
983	LexGLUE Paragraphs were split using double	1031
984	newlines as separators, preserving the inherent	1032
985	paragraph structure in the dataset.	1033
986	BookSum and SQUALITY Lines were joined	1034
987	with a blank space to create continuous text blocks.	1035
988	Sentences and paragraphs were split using the SaT-	
989	I3 model with a threshold of 0.8, producing lists	
990	of sentences grouped into paragraphs. Paragraphs	
991	exceeding 3,000 characters were further split us-	
992	ing the SaT-I3 model with a stricter threshold of	
993	0.4. Consecutive short paragraphs (fewer than 2	
994	sentences or under 400 characters) were merged	
995	to ensure coherence, especially for dialogue-heavy	
996	sections.	
997	JRC-Acquis Lines were joined with a blank	
998	space to preserve the flow of text. Sentences and	
999	paragraphs were split using the SaT-I3 model with a	
1000	standard threshold of 0.5. Consecutive paragraphs	
1001	containing fewer than 2 sentences were merged.	
1002	Sections and subsections were extracted from para-	
1003	graph beginnings using the dataset’s consistent	
1004	numbered format (e.g., 1.1.2), serving as structural	
1005	indicators.	
1006	MultiUN Lines were joined using blank spaces	
1007	to form initial text blocks. Sentences and para-	
1008	graphs were split using the SaT-I3 model with a	
1009	threshold of 0.5. Short consecutive paragraphs	
1010	(fewer than 2 sentences each, and up to 20 sen-	
1011	tences total) were merged to improve readability	
1012	and flow.	
1013	Wikipedia Original paragraphs were identified	
1014	using empty lines (meaning double newline in the	
1015	original text), which appeared as blank lines or	
1016	in the CSV format. Long paragraphs (over 500	
1017	tokens) were split further using the SaT-I3 model	
1018	with a threshold of 0.5 to improve segmentation	
1019	accuracy for longer text units.	
1020	GovReport Same as LexGLUE, paragraphs were	
1021	split using double newlines as separators.	
1022	Summscreen Initial paragraph segmentation was	
1023	based on scene indicators ([SCENE-BREAK]) in	
1024	the transcripts. However, this often resulted in ex-	
1025	cessively long paragraphs, with some documents	
	containing only one or two paragraphs. Text for-	1026
	matting issues, such as double spaces in punctua-	1027
	tion (e.g., " . "), were corrected to align with the	1028
	SaT-I3 model’s sensitivity. Long paragraphs ex-	1029
	ceeding 3,000 characters were re-segmented using	1030
	the SaT-I3 model with a threshold of 0.9. Short con-	1031
	secutive paragraphs containing only one sentence	1032
	were merged to form cohesive segments.	1033
	B Data Details	1034
	Our data collection	1035
	• 100% human annotated (no LLM pre-	1036
	selection)	1037
	• 7 domains (political, wikipedia, scientific, lit-	1038
	erature, conversational, legal	1039
	• 252 Source Documents (5k)	1040
	• 4 lengths of the same Source Document (\approx 5k,	1041
	\approx 1k, \approx 500, \approx 250 words)	1042
	• 13-15 QA on each Long Context Source Doc-	1043
	ument	1044
	• Annotation on the presence of these QA on	1045
	each of the summaries	1046
	• Human evaluation of automatic and manual	1047
	summaries	1048
	• Human evaluation of summary expansion	1049
	C Summary evaluation	1050
	Table 5 summarises the metrics used to evaluate.	1051
	D Summarization Guidelines	1052
	Annotator proficiency requirements	1053
	• Native speaker of English	1054
	• Editor / writer / domain expert	1055
	Task You will receive document(s) that are ap-	1056
	proximately 5,000 words or longer from the follow-	1057
	ing domains:	1058
	• Political (GovReports, MultiUN)	1059
	• News (Seahorse)	1060
	• Wikipedia (Wikipedia)	1061
	• Scientific/Technical (FacetSum)	1062
	• Literature (BookSum, SQuality)	1063

Task	Area	Metric	Description	Reference
Sum	Target similarity	R-L	ROUGE-L (longest common subsequence)	Lin (2004)
Sum/SumExp	Grammaticality	REP-3	Portion of duplicated N-grams (N=4)	Welleck et al. (2019)
Sum/SumExp	Fluency	CoLA	Sentence fluency classifier score	Krishna et al. (2020)
Sum/SumExp	Coherence	COH-2	2nd-order word-level coherence score	Parola et al. (2023)
Sum	Attribution	SH-4	Seahorse-Large-Q4 score	Clark et al. (2023)
Sum	Semantic coverage	SH-5	Seahorse-Large-Q5 coverage score	Clark et al. (2023)
SumExp	Work Count	WC		
Sum/SumExp	Overall	AVG	Empirical average of metrics	
Sum/SumExp	Overall	HE	HelloEval score	Que et al. (2024)

Table 5: Summary of automatic metrics used in different tasks.

- Conversational (Summscreen)

- Legal (LexGlue, JRCAcquis)

The documents will contain sections/chapters. You will need to summarize them retaining the section alignment. For certain domains, there will be additional guidance in the form of special guidelines (legal, medical etc.)

You will need to create 3 summaries:

- Summary 1: around 20% of the source text (1,000 words if total length is 5000)
- Summary 2: around 10% of the source text (500 words)
- Summary 3: around 5% of the source text (250 words)

After finishing summarizing, you will need to write a minimum of 15 questions with corresponding answers (QA) per each 5000 words.

Requirements for Summarization Here’s more information on what that means and how to summarize:

Please read the provided document in its entirety. Consider making notes of the main core ideas while you read. After you have finished reading, please write a short summary of the source text. The summary should:

- Be much shorter than the source text. Please see the information about the length above.
- Convey ALL the main core ideas and information of the source document.
- Have a structure of a standalone cohesive text.
- Follow the flow of the section/paragraph structure of the source text, try to summarize it from top to bottom.

- Each paragraph in each summary should be marked with a number of the source text section/chapter showing where this information is from.

Here’s a checklist which can help you with the task:

- **Understand the Main Idea:** Read the entire text to grasp the overall theme and the author’s intent. Identify the main idea of the text.
- **Highlight Key Points:** Mark or note down the essential points and arguments.
- **Eliminate Redundancies:** Remove any repetitive information or examples that do not add value to the understanding of the main idea.
- **Use Your Own Words:** Paraphrase the key points in your own words instead of copying verbatim. This helps ensure the summary is concise.
- **Keep It Objective:** Focus on the information presented in the text without inserting personal opinions or interpretations.
- **Structure the Summary:** Organize the ideas logically, maintaining the flow of the original text.
- **Be Concise:** Aim for clarity and brevity; Use simple and direct language to convey the points.
- **Review and Revise:** Compare the summary to the original text to ensure accuracy and completeness. Edit for coherence, transitions, and readability.

Additional Guidelines for Conversational Text

You may be assigned to work on conversational type of text, such as meeting transcripts, screenplays, and novels. Since the text structure is quite different from documents, this is an additional guideline to help you working on conversational text:

- Skim through the whole document: Try to get a rough idea of the whole plot
- Identify the characters and main core ideas: Identify the main characters and focus on their interaction
- Omit the trivial details: there maybe side plot or supporting characters in the source text, carefully decide if it is related to the main core ideas (plot)
- Group and summarize with respect to main core ideas: There could be plot twists or related hints in the source documents. Remember the summary should be clear and straightforward, the plot outline
- Should be clear in the summary without referencing to the source document

Requirements for Question and Answer sets

- The questions need to be abstractive not extractive: this means they need to be directed at the ideas in the text, not words and sentences as such.
- The questions should be open-ended, not Yes/No questions
- The correct answers should cover the main points in the source text: the questions should roughly correspond to paragraphs/sections in the text
- Thus, the correct answer should cover points reflected in the summary (for your convenience, you can refer to your longest summary, but please mind your shorter summary also need to be able to answer at least some of the questions)
- The answers should not be short (30 words or more)
- The correct answer should be found namely IN THE SUMMARY and not able to be just pulled from general knowledge

- If possible, refrain from factual questions, but try composing questions for reasoning, such as WHY- questions
- You are encouraged to combine information from different sections together
- Avoid only asking questions about the beginning and the end of the section, use the information in the middle as well
- If possible, the questions should not have several answer possibilities.

E Human Summary Sample

This is one sample of 3 human written summaries and QA set of the document 4586 from GovReport.

E.1 Source Document Excerpt

[This is the first and last paragraphs of the source document. The whole document is 5411 words long.]

This report examines technological innovation in payment systems generally and particular policy issues as a result of retail (i.e., point of sale) payment innovation. The report also discusses wholesale payment, clearing, and settlement systems that send payment messages between banks and transfer funds, including the "real-time payments" service being introduced by the Federal Reserve. This report includes an Appendix that describes interbank payment, clearing, and settlement systems related to U.S. payments.

...

To address systemic risk concerns, a private RTP system could be designated as a systemically important Financial Market Utility (FMU) under Title VIII of the Dodd-Frank Act (P.L. 111-203). The Dodd-Frank Act allows the Financial Stability Oversight Council , a council of financial regulators led by the Treasury Secretary, to designate a payment, clearing, or settlement system as systemically important on the grounds that "the failure of or a disruption to the functioning of the FMU could create or increase the risk of significant liquidity or credit problems spreading among financial institutions or markets and thereby threaten the stability of the U.S. financial system." FMUs, currently including the Clearing House Interbank Payments System, are subject to heightened regulation, and the Fed has supervisory and enforcement powers to ensure those standards are met. Policymakers could consider whether systemic risk concerns are

better addressed through Fed operation of payment and settlement systems or Fed regulation of private systems.

E.2 Long, Medium and Short Summaries

Each paragraph of the summaries is paired with the source paragraph id (e.g. p1, p2, etc) to indicate the information source.

Long Summary (20%) To the average consumer, swiping their credit card seems simple, because the complex the infrastructure involved in is 'hidden'. These deceptively "simple" electronic payments are comprised of three main steps. First, the sender makes the payment through an online payment service or an app, which instructs the sender's bank to make the payment to the recipient. Second, the bank sends a payment message to the recipient's bank through a payment system or clearing service. Finally, the payment is completed (settled) when the funds are received by the recipient. (p2)

Some of the bank-to-bank (ACH) payment, clearing, and settlement (PCS) systems are operated by the Federal Reserve, and others by private-sector organizations. Recently, the use of electronic payment methods (credit card, debit card, and ACH) has grown, while the use cash and check payments has declined. Electronic payments have been made easier and more convenient with digital wallets and payment apps like Venmo, Cash App, and Zelle - all of which require users to link a bank account, credit card, or debit card. (p4, p5, p6, p7)

There are concerns about whether current regulations are equipped to handle electronic payments. If not, this poses potential risks to cybersecurity, data privacy, industry competition, and consumer access and protection. Current payment regulations depend, in part, on if the service is provided by a bank, who have many strict regulatory requirements. As such, Nonbank payment systems are not subject to existing regulatory enforcement and can only be supervised - as money transmitters at state level and money service businesses at federal level. (p8, p9)

Electronic payment regulatory concerns could be addressed by including nonbank payment companies into the bank regulatory regime. One way could be via the Office of the Comptroller of the Currency (OCC) special purpose national bank charter. And another, through a state-level industrial loan company (ILC) charter with the Federal Deposit Insurance Corporation (FDIC). Both meth-

ods could provide nonbank firms access to the Fed wholesale payment systems, which could be advantageous. However, some state regulators have filed lawsuits to block nonbank companies access to these charters, arguing that it allows companies to circumvent state consumer protections. So far, no companies have applied for an OCC charter, likely due to the legal uncertainty surrounding it. (p11, p12)

The main argument against nonbank payment companies filing ILC charters is that it would allow them to own banks - and the FDIC has not approved deposit insurance for a new ILC since 2006. Opponents argue that allowing a company to own a bank could expose the US economy to risks like imprudent underwriting. Proponents assert that these concerns are exaggerated, noting that several other countries allow similar arrangement with no ill effects. So far, Square is the only company with a pending application and two other companies have withdrawn their applications. (p13)

New technology reduces some risks related to payments but creates new ones. The risk of having one's wallet stolen is reduced, but payment information is subject to more sophisticated risks such as malware attacks. Furthermore, storing payment information on a variety of websites, apps, and devices creates more opportunities for hackers. After recent security breaches which allowed user information to be stolen, several solutions have been proposed. For example, a federal breach notification law could be enacted, to create federal cybersecurity standards or to increase penalties for companies with inadequate security measures. (p15, p16, p17)

Payment systems need to collect detailed information about customer transactions in order to function properly. This data can be used by companies to target ads. Scammers can also use this information for fraudulent purposes. The constantly increasing use of Electronic payments has led to questions about how user data is used and whether consumers are sufficiently informed and given enough control about how their data is used. (p18)

There are some consumer benefits to storing consumer data. It can help them track payments and budget more easily by importing to budgeting apps. They can also share financial information with banks more easily when applying for loans. But, given the benefits and the risks, the question remains: how much access should companies have to individuals' information? (p19)

Privacy policies are another area of concern with respect to consumer protection and electronic payments. According to the Bureau of Consumer Financial Protection (CFPB), it is difficult to provide disclosures that are clear and easy to understand, partly due to the small screens on phones. Clearer privacy policies and allowing consumers more control over how their data is used could help. (p20)

The Electronic Fund Transfer Act, Regulation E implemented by the CFPB, is the most relevant law aimed at protecting consumers who are making electronic payments. It mandates consumer disclosures, limits consumer liability for unauthorized payments, and maintains procedures for resolving errors. Further regulations are being considered. (p22)

Consumers could also be protected through financial education, especially for more at-risk older and lower-income groups. This could include learning how to use new payment systems safely and how to protect against financial harm, as well as knowing how to get help if something goes wrong. (p24)

Payment system innovations may affect consumers differently based on income. Consumers who mainly pay with cash, don't have bank accounts, or don't have internet or mobile access won't be able to benefit. Neither will those who are not comfortable using new technology. (p25)

However, surveys reveal that 83% of underbanked, and 50% of unbanked, consumers have smart phone access. So, as costs of these payment services decline, some marginalized groups could experience better access to the the financial system through access to digital currency channels via cash equivalents like pre-paid cards. But, the cost of internet and mobile data plans may limit access to faster payment systems, so this also needs to be considered. (p26, p27, p28)

Faster payment systems may also benefit low-income consumers by allowing them faster access to their paychecks and other fund transfers. But a potential drawback is that withdrawals from their accounts would occur more quickly as well. (p28)

In 2019, the Fed announced that it plans to create a wholesale real-time payment (RTP) system. (p32)

Originally, the Fed's primary function was to provide bank-to-bank check-clearing services. Private clearing houses were experiencing issues that led to the creation of the Fed. As payment methods have evolved, the Fed has begun providing other types of payment systems. It does this by linking

the accounts that all banks keep at the Fed so that it can complete the transfers. The new system that the Fed is developing, called FedNow, would allow payments to occur in real time, rather than later in the day - or even the next day, as is the case currently. (p33, p35)

However, there are some concerns regarding implementation of FedNow. Many worry that it will undermine private sector development of similar systems. Others fear that failing to implement FedNow will lead to a monopoly of a private-sector company, to the detriment of consumers and smaller banks. (p42, p44)

Medium Summary (10%) Electronic payments have three stages. First, the sender makes the payment through an online payment service or an app, which instructs the sender's bank to make the payment to the recipient. Second, the bank sends a payment message to the recipient's bank through a payment system or clearing service. Finally, the payment is completed (settled) when the funds are received by the recipient. (p2)

Some of the bank-to-bank (ACH) payment, clearing, and settlement (PCS) systems are operated by the Federal Reserve, and others by private-sector organizations. Recently, the use of electronic payment methods (credit card, debit card, and ACH) has grown, while the use cash and check payments has declined. Electronic payments have been made easier and more convenient with digital wallets and payment apps like Venmo, Cash App, and Zelle - all of which require users to link a bank account, credit card, or debit card. (p4, p5, p6, p7)

There is concern about whether current regulations are equipped to handle these technological advances. If not, they could pose risks to cybersecurity, data privacy, industry competition, and consumer access and protection. (p8)

One way to address these concerns is to add nonbank companies to the bank regulatory regime. Another is via the Office of the Comptroller of the Currency (OCC) special purpose national bank charter. And another, through a state-level industrial loan company (ILC) charter with the Federal Deposit Insurance Corporation (FDIC). Both could provide nonbank firms access to the Fed wholesale payment systems, which could be advantageous. However, some state regulators have tried to block nonbank access to these charters, arguing that it allows companies to circumvent state consumer protections. So far, no companies have applied for

1426	an OCC charter, likely due to the legal uncertainty	1477
1427	surrounding it. (p11, p12)	1478
1428	The main argument against ILC charters is that	1479
1429	it would allow companies to own banks. The FDIC	1480
1430	has not approved deposit insurance for a new ILC	1481
1431	since 2006. Opponents argue that allowing a re-	1482
1432	tailer to own a bank could expose the US economy	1483
1433	to risks such as imprudent underwriting. Propo-	1484
1434	nents assert that these concerns are exaggerated	1485
1435	and that several other countries allow similar ar-	1486
1436	rangement with no ill effects. Currently, Square	1487
1437	is the only company with a pending application.	1488
1438	(p13)	1489
1439	Privacy policies are another area of concern. Ac-	1490
1440	cording to the Bureau of Consumer Financial Pro-	1491
1441	tection (CFPB), it is difficult to provide disclosures	1492
1442	that are clear and easy to understand, partly due to	1493
1443	the small screens on phones. Clear privacy policies	1494
1444	and allowing consumers more control over how	1495
1445	their data is used could help. (p20)	1496
1446	The Electronic Fund Transfer Act, Regulation	1497
1447	E implemented by the CFPB, is the most relevant	1498
1448	law aimed at protecting consumers who are making	1499
1449	electronic payments. It mandates consumer disclo-	1500
1450	tures, limits consumer liability for unauthorized	1501
1451	payments, and maintains procedures for resolving	1502
1452	errors. (p22)	1503
1453	In 2019, the Fed announced that it plans to create	1504
1454	a wholesale real-time payment (RTP) system. The	1505
1455	proposed system, called FedNow, would allow pay-	1506
1456	ments to occur in real time, rather than later in the	
1457	day - or even the next day, as is the case currently.	1507
1458	(p32, p33, p35)	
1459	There are some concerns regarding FedNow.	
1460	Many worry that it will undermine private sector de-	
1461	velopment of similar systems. Others fear that fail-	
1462	ing to implement FedNow will lead to a monopoly	
1463	of a private-sector company, to the detriment of	
1464	consumers and smaller banks. (p42, p44)	
1465	Short Summaries (5%) Electronic payments	
1466	have three stages. First, the sender makes the pay-	
1467	ment through an online payment service or an app,	
1468	which instructs the sender's bank to make the pay-	
1469	ment to the recipient. Second, the bank sends a	
1470	payment message to the recipient's bank through	
1471	a payment system or clearing service. Finally, the	
1472	payment is completed (settled) when the funds are	
1473	received by the recipient. (p2)	
1474	Some of the bank-to-bank payment, clearing,	
1475	and settlement (PCS) systems are operated by the	
1476	Federal Reserve, and others by private-sector orga-	
	nizations. Recently, the use of electronic payment	1477
	methods has grown, while the use cash and check	1478
	payments has declined. (p4)	1479
	There is concern about whether current regula-	1480
	tions are equipped to handle these technological	1481
	advances. If not, they could pose risks to cyber-	1482
	security, data privacy, industry competition, and	1483
	consumer access and protection. (p8, p9)	1484
	A potential way to resolve concerns regarding	1485
	regulation is to add certain nonbank payment com-	1486
	panies into the bank regulatory regime. One way	1487
	to accomplish this is through the OCC or Office of	1488
	the Comptroller of the Currency special purpose	1489
	national bank charter. The second is through a state-	1490
	level industrial loan company (ILC) charter with	1491
	the Federal Deposit Insurance Corporation (FDIC).	1492
	(p11, p12)	1493
	Privacy policies are another area of concern. Ac-	1494
	cording to the Bureau of Consumer Financial Pro-	1495
	tection (CFPB), it is difficult to provide disclosures	1496
	that are clear and easy to understand, partly due to	1497
	the small screens on phones. Clear privacy policies	1498
	and allowing consumers more control over how	1499
	their data is used could help. (p20)	1500
	In 2019, the Fed announced that it plans to create	1501
	a wholesale real-time payment (RTP) system. The	1502
	proposed system, called FedNow, would allow pay-	1503
	ments to occur in real time, rather than later in the	1504
	day - or even the next day as is the case currently.	1505
	(p32, p33, p35)	1506
	E.3 Question and Answer Set	1507
	Question 1: What are the three parts of a payment	1508
	system?	1509
	• Answer: First, there is the sender or the per-	1510
	son making the payment through an online	1511
	payment service or an app, which instructs	1512
	the sender's bank to make the payment to	1513
	the recipient. Second, the bank sends a pay-	1514
	ment message to the recipient's bank through	1515
	a payment system or clearing service. Finally,	1516
	the payment is completed when the funds are	1517
	transferred, or settled.	1518
	• (Information contained in 20% Summary,	1519
	10% Summary, 5% Summary.)	1520
	• (Source paragraph number: p2)	1521
	Question 2: Who operates bank-to-bank payment,	1522
	clearing, and settlement systems?	1523

1524	• Answer: Some of these systems are operated	• (Information contained in 20% Summary,	1568
1525	by the Federal Reserve and some are operated	10% Summary, 5% Summary.)	1569
1526	by private-sector organizations.		
1527	• (Information contained in 20% Summary,	• (Source paragraph number: p29)	1570
1528	10% Summary, 5% Summary.)		
1529	• (Source paragraph number: p4)	Question 7: What are some reasons for the increase	1571
1530	Question 3: What issues could there be if current	in electronic payments?	1572
1531	regulations are not equipped to handle these pay-		
1532	ment system innovations?	• Answer: Electronic payments have increased	1573
1533	• Answer: If regulations are inadequate, there	because of payment apps such as Venmo,	1574
1534	could be issues related to cybersecurity, data	Cash App, and Zelle make it convenient and	1575
1535	privacy, industry competition, and consumer	easy for consumers to send payments. Digital	1576
1536	access and protection.	wallets stored on phones are another reason	1577
1537	• (Information contained in 20% Summary,	for increased electronic payments due their	1578
1538	10% Summary, 5% Summary.)	ease of use and convenience.	1579
1539	• (Source paragraph number: p8)	• (Information contained in 20% Summary,	1580
1540	Question 4: What are two ways to bring nonbank	10% Summary.)	1581
1541	companies into the bank regulatory regime?	• (Source paragraph number: p6)	1582
1542	• Answer: One way to accomplish this is	Question 8: What is necessary in order for a con-	1583
1543	through the OCC or Office of the Comptroller	sumer to be able to use electronic payment ser-	1584
1544	of the Currency special purpose national bank	vices?	1585
1545	charter. The second is through a state-level	• Answer: The consumer must have a debit	1586
1546	industrial loan company (ILC) charter with	card, credit card, or bank account linked to an	1587
1547	the Federal Deposit Insurance Corporation	electronic payment system.	1588
1548	(FDIC).	• (Information contained in 20% Summary,	1589
1549	• (Information contained in 20% Summary,	10% Summary.)	1590
1550	10% Summary, 5% Summary.)	• (Source paragraph number: p7)	1591
1551	• (Source paragraph number: p11,12)	Question 9: Why have state regulators filed law-	1592
1552	Question 5: According to the Bureau of Consumer	suits to block the OCC?	1593
1553	Financial Protection, what are some of the difficul-	• Answer: Regulators feel that the OCC charter	1594
1554	ties with privacy policies?	would allow companies to avoid state regula-	1595
1555	• Answer: It is difficult to provide disclosures	tions that protect consumers.	1596
1556	that are clear and easy to understand, partly	• (Information contained in 20% Summary,	1597
1557	due to the small screens on phones.	10% Summary.)	1598
1558	• (Information contained in 20% Summary,	• (Source paragraph number: p11, p12)	1599
1559	10% Summary, 5% Summary.)	Question 10: What is the main argument against	1600
1560	• (Source paragraph number: p20)	the ILC charter?	1601
1561	Question 6: What is FedNow?	• Answer: The ILC would allow companies	1602
1562	• Answer: In 2019, the Fed announced that it	such as retailers to own banks. Opponents are	1603
1563	plans to create a wholesale real-time payment	concerned that this could lead to imprudent	1604
1564	(RTP) system. The proposed system, called	underwriting and could hurt the US economy	1605
1565	FedNow, would allow payments to occur in	by exposing it to risk.	1606
1566	real time, rather than later in the day or even	• (Information contained in 20% Summary,	1607
1567	the next day as is the case currently.	10% Summary.)	1608

1609	• (Source paragraph number: p13)	• (Information contained in 20% Summary.)	1651
1610	Question 11: What does the Electronic Funds	• (Source paragraph number: p28)	1652
1611	Transfer Act Regulation E do?		
1612	• Answer: Regulation E mandates consumer	F Summarization Human Evaluation	1653
1613	disclosures, limits consumer liability for unau-	Guidelines	1654
1614	thorized payments, and maintains procedures	Annotator proficiency requirements All anno-	1655
1615	for resolving errors.	tators must meet ALL of the following require-	1656
1616	• (Information contained in 20% Summary,	ments:	1657
1617	10% Summary.)	• Native speaker of English AND	1658
1618	• (Source paragraph number: p22)	• Language related degree holder or related pro-	1659
1619	Question 12: How could financial education help	fessionals	1660
1620	consumers use electronic payment systems safely?	Background Information What is a good sum-	1661
1621	• Answer: Consumers could be taught how to	mary? A good summary should meet the following	1662
1622	use new payment systems safely and how to	criteria:	1663
1623	protect against financial harm, as well as how	• Conciseness: The summary should only con-	1664
1624	to get help if something goes wrong.	tain the most important information while	1665
1625	• (Information contained in 20% Summary.)	maintaining readability. Trivial information	1666
1626	• (Source paragraph number: p24)	should not be included, even in the longest	1667
1627	Question 13: What are is an argument against the	summary. Additionally, the summary should	1668
1628	FedNow?	be comprehensible on its own, without need-	1669
1629	• Answer: Many worry that it will undermine	ing to refer to additional documentation.	1670
1630	private sector development of similar systems.	• Coverage of main core ideas: The summary	1671
1631	• (Information contained in 20% Summary,	should preserve the most important ideas, re-	1672
1632	10% Summary.)	gardless of its length. In our task, summaries	1673
1633	• (Source paragraph number: p42, p44)	are created by gradually omitting less impor-	1674
1634	Question 14: How can storing more consumer data	tant information. Therefore, we expect that	1675
1635	benefit consumers?	the core ideas will be retained in all sum-	1676
1636	• Answer: It can help consumers track pay-	maries, even the shortest ones. Main core	1677
1637	ments and budget more easily using budgeting	ideas should be the key ideas that help the	1678
1638	apps. They can also share financial informa-	reader to understand the main topic. Depend-	1679
1639	tion with banks more easily when applying	ing on the type of documents, the definition	1680
1640	for loans.	of idea would be slightly different. For ex-	1681
1641	• (Information contained in 20% Summary.)	ample, if the source document is a meeting	1682
1642	• (Source paragraph number: p19)	note, the summary should include the main	1683
1643	Question 15: How could faster payment systems	topic, the discussion, the result / final decision.	1684
1644	affect low-income consumers?	The trivial details like greetings or small talks	1685
1645	• Answer: Faster payment systems may bene-	should not be included. If the source docu-	1686
1646	fit low-income consumers by allowing them	ment is a novel, the summary should focus on	1687
1647	faster access to their paychecks and other fund	main characters and important events rather	1688
1648	transfers. But a potential drawback is that	than trivial description of the character or side	1689
1649	withdrawals from their accounts would occur	events.	1690
1650	more quickly as well.	• Attribution: the information in the summary	1691
		can be accurately referred back to the source	1692
		documents. All the information in the sum-	1693
		mary should be an abstraction from the source	1694
		documents. No additional information that	1695
		can not be found in the source document	1696
		should be included in the summary.	1697

- **Cohesion as a document:** Each summary should be an abstraction of the entire source document. All the information or ideas should be digested from different parts of the source document and combined into a new paragraph. Merely shortening a document paragraph by paragraph will not be considered as a good summary. Similarly, a bulletin-like document jumping from point to point also will not be considered as a good summary.

Annotation Our summarization structure is as follows: The source text which is approximately 5,000 words long gets summarized three times: The first summary is 20% of the original length of the doc. It should retain all the core ideas of the source document. The second summary is 10% of the length of the source document. It should also retain all the core ideas of the source. The third summary is short, it should be 5% of the source length. We understand there will be some information loss, but again, all the core ideas should be present in the summary.

There are two tasks related to evaluating the summary.

Task 1 In this task, you need to rate the overall quality of the summaries regarding several aspects. Here is the detailed workflow:

Step 1 Screening

Please spend no more than 5 minutes skimming through the longest (20%) summary and answer the question below.

- Q1: Is it a cohesive text? Can you fully understand it?
 - If NOT, reject the task completely.
 - If YES, continue with the following steps

Step 2 Read the texts and take notes

Please read the whole source document carefully and take notes in your own way. It could be highlighting the key points or jotting down the ideas in your own words or any means that can help you digest the document. While reading please do not skip any line. After reading and taking note, you should be able to identify several main core ideas or more (You can spot more main core ideas if the text is longer). Please continue to read the summary and identify if the main ideas also exist in the summaries. Now check how many ideas can be found in the summary.

(This procedure is to help rate the summaries more objectively, you are not required to submit the highlights or the notes.)

Step 3 Rate the summaries

Answer all the following questions with a 4-point scale:

- Q2a Check the attribution of the summary. Can all the information in the summary be attributed to the source text?
 - Give 4 points if yes, all the information can be directly attributed to the source text.
 - Give 3 points if mostly yes, only 1 idea seems to not be found in the source text.
 - Give 2 points if not really, more than 1 idea cannot be attributed to the source text.
 - Give 1 point if not, most ideas cannot be found in the source text and seem to be completely new.
 - Give 0 points if not, none of the ideas can be found in the source text.
- Q2b Check the coverage of main core ideas of the source text in the summary. Are all the main core ideas of the source document retained?
 - Give 4 points if yes, all the main core ideas of the source are retained.
 - Give 3 points if mostly yes, only 1 or 2 main core ideas are not found in the summary.
 - Give 2 points if not really, more than 2 main core ideas are not found in the summary.
 - Give 1 point if not, most main core ideas are not found in the summary.
 - Give 0 points if not, none of the ideas can be found in the summary
- Q2c Check the conciseness of the summary. Is the summary short and clear without repetition and redundancy?
 - Give 4 points if yes, the summary is not wordy but clear.
 - Give 3 points if mostly yes, but 1 part is unnecessary.
 - Give 2 points if not really, more than 1 part is unnecessary.

1885	as an interesting plot. For example, you can	– Give 4 points if yes, it is a well structure	1932
1886	check the following questions depending on	screenplay / novel	1933
1887	the story: If it's a comedy, does it sound funny	– Give 3 points if mostly yes, but 1 part is	1934
1888	to you? If it's a romance story, does it evoke	missing.	1935
1889	the proper sentiment? Are the added details	– Give 2 points if not really, more than 1	1936
1890	aligned with the plot, or do they feel out of	part is missing.	1937
1891	place? ... etc	– Give 1 point if text does not follow the	1938
1892		structure of a screenplay / novel	1939
1893	Task 1 In this task, you need to rate the over-	– Give 0 points if not, the text doesn't not	1940
1894	all quality of the summary expansions regarding	read as a cohesive text at all	1941
1895	several aspects. Here is the detailed workflow:		
1896	Step 1 Screening	• Q2c Check the richness in details. Does it	1942
1897	Please spend no more than 5 minutes skimming	contain enough details?	1943
1898	through the long form text and answer the question		
1899	below.	– Give 4 points if yes, the text contains a	1944
1900	• Q1: Is it a cohesive text? Can you fully under-	lot of details.	1945
1901	stand it?	– Give 3 points if yes, the text contains	1946
1902	– If NOT, reject the task completely.	details but has room for improvement.	1947
1903	– If YES, continue with the following steps	– Give 2 points if not really, the text con-	1948
1904		tains limited details.	1949
1905	Step 2 Read the texts and highlight key points	– Give 1 point if not, the text contains very	1950
1906	Please read the original summary carefully and	few details	1951
1907	highlight the key points and make notes. Do not	– Give 0 points if not, the text does not	1952
1908	skip any line. Continue to read other summaries	provide any additional details at all.	1953
1909	and long form text, highlighting the key points that		
1910	are the same as the original summary	• Q2d Check the creativity. Does the added	1954
1911	(This procedure is to help rating the summaries	details novel and original while being relevant	1955
1912	more objectively, you are not required to submit	to the core main ideas?	1956
1913	the highlights or the notes.)		
1914	Step 3 Rate the long form text	– Give 4 points if yes, all of the added de-	1957
1915	Answer all the following questions with a 4-	tails are novel and original	1958
1916	point scale separately for the long form text:	– Give 3 points if yes, most of the added	1959
1917		details are novel and original.	1960
1918	• Q2a Check the coverage of main core ideas.	– Give 2 points if not really, only some of	1961
1919	Are all the core concepts of the original sum-	the added details are novel and original.	1962
1920	mary retained?	– Give 1 point if no, very few added details	1963
1921	– Give 4 points if yes, all the main core	are repetitive.	1964
1922	ideas are retained.	– Give 0 points if no, no added details are	1965
1923	– Give 3 points if mostly yes, only 1 or 2	novel and original .	1966
1924	core ideas are lost.		
1925	– Give 2 points if not really, more than 2	• Q2e Check the non-repetitiveness. Does it	1967
1926	ideas are lost.	repeat a lot?	1968
1927	– Give 1 point if not, most ideas are lost.		
1928	– Give 0 points if not, none of the ideas	– Give 4 points if no, all of the details are	1969
1929	can be found in the source text.	unique and different from each other	1970
1930		– Give 3 points if no, most of the details	1971
1931	• Q2b Check the cohesion of text. Is it well	are unique, only one is repeated	1972
	structured? Does it contain all the necessary	– Give 2 points if not really, some of the	1973
	components? (Scene description, dialog, main	details are repetitive	1974
	characters, etc) Does it flow logically and	– Give 1 point if yes, most the added de-	1975
	maintain consistency?	tails are repetitive	1976

1977		
1978		– Give 0 points if not, all the added details are repetitive
1979	• Q2f Rate the story plot. How interesting is it to you? Is it engaging and compelling?	
1980		
1981		– Give 4 points if the text is very interesting.
1982		
1983		– Give 3 points if the text is quite interesting.
1984		
1985		– Give 2 points if the text is somewhat interesting.
1986		
1987		– Give 1 point if the text is only slightly interesting.
1988		
1989		– Give 0 points if the text is dull and not interesting at all.
1990		
1991	After evaluating all the aspects of the expanded text, please give an overall score of 0-10 on the quality of the expanded text.	
1992		
1993		
1994	• Q3 Do you think expanded text is well written? Do you think it is a good read? On a scale of 0- 10, how would you rate the overall quality?	
1995		
1996		
1997		
1998		– 10: The text is perfect in every aspects
1999		– 8-9: The text is considered good. It contains minor issues in certain aspects but it meets all requirements with room for improvement.
2000		
2001		– 6-7: The text is moderate, it contains non-critical errors but to help the reader understand the source documents
2002		
2003		– 4-5: The text is below acceptable level. It contains critical errors that cause trouble to read
2004		
2005		– 2-3: The text contains very limited information that is relevant to the source summary
2006		
2007		– 0-1: The text is barely readable and comprehensible or it barely contains relevant information to the source summary.
2008		
2009		
2010		
2011		
2012		
2013		
2014		
2015	Make sure you have answered all the questions for every expanded summary and long form text.	
2016		
2017	Task 2 You will be provided with around 15 questions depending on the length of the documents. You need to answer YES or NO to each QA set. Answer YES only when the answer is directly stated, heavily implied, or logically entailed in the text.	
2018		
2019		
2020		
2021		

H Prompting Details

The prompts contain three parts: General guideline, domain-specific prompts, and input context. The general guideline adapts the human guidelines (Appendix D) for the summarization and summary expansion, while the domain specific prompts give extra information about the domain as instructions of expected output. In the prompt template below, the general guideline is provided, `{{domain-X}}` denotes the domain-specific prompt. `{{input}}` is for the input document for the summarization task and human summaries for the summary expansion task.

Prompt for summarization

```

"""
You are a professional editor and reader.
You are reading a {{domain}}
{{domain-meta}}.
The {{domain}} starts with [START]
and ends with [END].
After you have finished reading, please
provide a summary of the {{domain}}.
{{domain-expect}}.
Make sure the summary has
{{len(input) * ratio + 200}} words or
less.
[START]
{{input}}
[END].
Write at least {{len(input) * ratio}}
words.
"""

```

Prompt for summary expansion

```

"""
You are a professional editor and reader.
You are reading a summary of {{domain}}
{{domain-meta}}.
The {{domain}} starts with [START]
and ends with [END].
After you have finished reading, write
a well-structured, consistent {{domain}}
that extends the summary.
{{domain-expect-expand}}.
[START]
{{input}}
[END].
Write at least {{len(source)}} words.
"""

```

The model-specific prompts for each domain are listed below. Note that not all domains have the prompt template for summary expansion.

• BookSum:

```

{{domain}}: "book chapter"
{{domain-meta}}: """about the book
[BOOK-TITLE], chapter [CHAP-NO],
title [[CHAP-TITLE]]."""
{{domain-expect}}: ""
{{domain-expect-expand}}: """Please
keep the main plot and characters

```

2081 if found in the summary.
2082 """.

2083 • **LexGLUE:**

2084 {{domain}}: "legal document"
2085 {{domain-meta}}: ""
2086 {{domain-expect}}: ""Keep the main
2087 ideas and terms in the document.
2088 """

2089 • **SQuALITY:**

2090 {{domain}}: "short story"
2091 {{domain-meta}}: ""
2092 {{domain-expect}} : ""Keep the main
2093 character names and narratives
2094 of the story.
2095 """
2096 {{domain-expect-expand}}: (same)

2097 • **Seahorse:**

2098 {{domain}}: "news article"
2099 {{domain-meta}}: ""
2100 {{domain-expect}}: ""

2101 • **FacetSum:**

2102 {{domain}}: "academic article"
2103 {{domain-meta}}: "about [TITLE]"
2104 {{domain-expect}}: ""Keep the
2105 structure of sections [SECTIONS]
2106 """
2107 {{domain-expect-expand}}: (same)

2108 • **JRC-Acquis:**

2109 {{domain}}: "document"
2110 {{domain-meta}}: "from European Commision"
2111 {{domain-expect}}: ""

2112 • **MultiUN:**

2113 {{domain}}: "document"
2114 {{domain-meta}}: "from United Nation"
2115 {{domain-expect}}: ""

2116 • **GovReport:**

2117 {{domain}}: "government report"
2118 {{domain-meta}}: ""
2119 {{domain-expect}}: ""

2120 • **Wikipedia:**

2121 {{domain}}: "Wikipedia article"
2122 {{domain-meta}}: "about [TITLE]"
2123 {{domain-expect}}: ""

2124 • **Summscreen:**

2125 {{domain}}: "screenplay"
2126 {{domain-meta}}: "about [TITLE]"
2127 {{domain-expect}}: ""Keep the main
2128 plot and characters in the screenplay
2129 """
2130 {{domain-expect-expand}}: ""Keep the
2131 main plot and characters in the
2132 summary.
2133 Write in the dialogue form with
2134 multiple utterances
2135 """

I Detailed Results

Table 6 shows the detailed results of different models in the summary expansion task. It is shown that, despite having repeated instructions on the length, all models behave greatly differently in different domains. In particular, GPT-4o-mini can generate longer scientific/technical texts, but struggle to generate longer texts in conversational texts without sacrificing the qualities. On the other hand, medium-sized models such as Llama-3.1-8B generate more texts consistently across domains, but at a higher repetition.

The results of the summarization task of different levels are detailed in tables 7, 8 and 9. According to human evaluation, the best performing result in Table 7 is with GPT-4o-mini in the wikipedia domain and in Table 8 is with the same model in the legal domain (LexGlue). For table 9, the best results are with human output in the conversational domain (Summscreen).

DATASET	Model	% WC	REP-3(↓)	CoLA↑	COH-2↑	AVG↑	HE↑	HUM↑
BookSum	GPT-4o-mini	0.641	0.459	0.960	0.739	0.536	87.161	6.691
	Llama-3.1-70B	1.539	0.650	0.867	0.842	0.526	50.969	5.123
	Llama-3.1-8B	1.695	0.736	0.936	0.861	0.550	53.692	5.160
SQuALITY	GPT-4o-mini	3.014	0.513	0.961	0.738	0.532	60.385	6.320
	Llama-3.1-70B	1.107	0.582	0.952	0.780	0.539	37.775	5.434
	Llama-3.1-8B	1.351	0.735	0.955	0.817	0.542	35.017	5.360
FacetSum	GPT-4o-mini	0.502	0.642	0.954	0.640	0.488	70.061	7.693
	Llama-3.1-70B	0.138	0.705	0.874	0.877	0.537	26.587	3.453
	Llama-3.1-8B	0.205	0.935	0.871	0.871	0.518	28.462	4.173
Summscreen	GPT-4o-mini	3.569	1.215	0.776	0.319	0.284	65.977	5.000
	Llama-3.1-70B	1.449	0.782	0.814	0.500	0.386	41.466	3.800
	Llama-3.1-8B	1.495	0.828	0.851	0.568	0.418	36.493	4.480

Table 6: Performance on the summary expansion task per dataset.

DATASET	Model	R-L(↑)	REP-3(↓)	CoLA↑	COH-2↑	SH-4↑	SH-5↑	AVG↑	HE↑	HUM↑
LCFO.5%										
LexGLUE	Human	n/a	0.258	0.930	0.807	0.617	0.339	0.528	46.690	6.360
	GPT-4o-mini	0.342	0.407	0.956	0.688	0.657	0.500	0.479	78.916	7.747
	Llama-3.1-70B	0.386	0.415	0.954	0.875	0.625	0.369	0.482	63.280	6.987
	Llama-3.1-8B	0.378	0.471	0.972	0.879	0.617	0.383	0.476	59.455	6.907
BookSum	Human	n/a	0.226	0.913	0.762	0.572	0.315	0.503	71.006	6.691
	GPT-4o-mini	0.302	0.257	0.977	0.857	0.599	0.485	0.532	93.168	6.815
	Llama-3.1-70B	0.377	0.362	0.976	0.846	0.578	0.374	0.483	76.871	6.272
	Llama-3.1-8B	0.372	0.400	0.973	0.846	0.581	0.347	0.469	72.999	6.049
SQuALITY	Human	n/a	0.263	0.922	0.760	0.520	0.334	0.497	33.534	5.173
	GPT-4o-mini	0.285	0.284	0.980	0.841	0.548	0.375	0.492	74.618	6.600
	Llama-3.1-70B	0.340	0.472	0.961	0.802	0.463	0.201	0.391	64.237	5.227
	Llama-3.1-8B	0.339	0.535	0.968	0.819	0.488	0.233	0.395	57.288	5.827
FacetSum	Human	n/a	0.260	0.945	0.835	0.691	0.436	0.571	57.456	7.053
	GPT-4o-mini	0.404	0.354	0.921	0.568	0.682	0.524	0.468	73.968	7.434
	Llama-3.1-70B	0.412	0.387	0.962	0.884	0.696	0.508	0.533	67.585	6.213
	Llama-3.1-8B	0.419	0.425	0.967	0.888	0.704	0.518	0.530	69.176	6.733
JRC-Acquis	Human	n/a	0.247	0.949	0.849	0.672	0.464	0.577	52.092	7.180
	GPT-4o-mini	0.352	0.383	0.952	0.539	0.682	0.593	0.477	82.239	7.347
	Llama-3.1-70B	0.390	0.424	0.942	0.883	0.690	0.470	0.512	60.948	6.306
	Llama-3.1-8B	0.368	0.427	0.945	0.882	0.673	0.449	0.504	59.209	6.514
MultiUN	Human	n/a	0.255	0.927	0.862	0.592	0.276	0.521	44.466	6.861
	GPT-4o-mini	0.352	0.364	0.968	0.549	0.630	0.528	0.462	76.639	7.347
	Llama-3.1-70B	0.402	0.400	0.955	0.903	0.618	0.303	0.476	76.121	6.611
	Llama-3.1-8B	0.378	0.443	0.965	0.907	0.608	0.320	0.471	59.683	6.806
Wikipedia	Human	n/a	0.246	0.961	0.810	0.664	0.246	0.527	68.484	6.893
	GPT-4o-mini	0.341	0.332	0.974	0.756	0.693	0.423	0.503	80.633	7.754
	Llama-3.1-70B	0.382	0.405	0.968	0.821	0.660	0.299	0.469	59.334	6.551
	Llama-3.1-8B	0.379	0.439	0.963	0.839	0.672	0.282	0.463	59.259	5.841
GovReport	Human	n/a	0.226	0.958	0.803	0.639	0.336	0.538	35.157	6.720
	GPT-4o-mini	0.340	0.333	0.978	0.722	0.696	0.538	0.520	81.420	7.280
	Llama-3.1-70B	0.407	0.363	0.973	0.870	0.651	0.430	0.512	54.626	6.080
	Llama-3.1-8B	0.407	0.353	0.971	0.855	0.620	0.354	0.489	52.231	6.44
Summscreen	Human	n/a	0.243	0.927	0.739	0.532	0.384	0.507	62.003	7.040
	GPT-4o-mini	0.294	0.289	0.984	0.832	0.514	0.346	0.478	60.347	6.627
	Llama-3.1-70B	0.390	0.328	0.985	0.849	0.523	0.259	0.458	70.638	6.173
	Llama-3.1-8B	0.375	0.373	0.976	0.854	0.526	0.266	0.450	69.116	5.667
Seahorse	Human	n/a	0.213	0.950	0.819	0.651	0.440	0.563	51.057	6.200
	GPT-4o-mini	0.295	0.279	0.985	0.832	0.647	0.556	0.548	67.220	7.613
	Llama-3.1-70B	0.352	0.382	0.965	0.842	0.661	0.434	0.504	65.295	6.293
	Llama-3.1-8B	0.354	0.369	0.978	0.846	0.649	0.472	0.515	65.893	6.427

Table 7: Performance on the 5% summarization task per dataset.

DATASET	Model	R-L(↑)	REP-3(↓)	CoLA↑	COH-2↑	SH-4↑	SH-5↑	AVG↑	HE↑	HUM↑
LCFO.10%										
LexGLUE	Human	n/a	0.351	0.940	0.829	0.660	0.362	0.544	62.141	7.387
	GPT-4o-mini	0.419	0.494	0.947	0.599	0.633	0.500	0.437	80.094	8.120
	Llama-3.1-70B	0.452	0.566	0.942	0.882	0.625	0.397	0.456	59.666	7.080
	Llama-3.1-8B	0.439	0.641	0.967	0.876	0.621	0.376	0.440	59.349	7.200
BookSum	Human	n/a	0.278	0.907	0.757	0.610	0.342	0.512	83.776	7.601
	GPT-4o-mini	0.327	0.308	0.978	0.858	0.578	0.442	0.510	94.867	7.062
	Llama-3.1-70B	0.427	0.456	0.967	0.835	0.573	0.335	0.451	82.114	6.469
	Llama-3.1-8B	0.415	0.511	0.966	0.844	0.551	0.313	0.432	73.681	6.420
SQuALITY	Human	n/a	0.313	0.917	0.773	0.548	0.312	0.497	51.501	6.000
	GPT-4o-mini	0.327	0.329	0.975	0.819	0.525	0.341	0.466	76.441	6.467
	Llama-3.1-70B	0.382	0.367	0.974	0.836	0.518	0.320	0.456	46.731	4.613
	Llama-3.1-8B	0.373	0.406	0.979	0.856	0.526	0.324	0.456	50.129	5.280
FacetSum	Human	n/a	0.328	0.942	0.840	0.710	0.425	0.570	69.570	7.680
	GPT-4o-mini	0.461	0.409	0.945	0.658	0.666	0.506	0.473	78.381	7.882
	Llama-3.1-70B	0.455	0.538	0.934	0.890	0.696	0.527	0.502	63.663	6.501
	Llama-3.1-8B	0.449	0.547	0.954	0.892	0.698	0.496	0.499	63.768	6.987
JRC-Acquis	Human	n/a	0.339	0.959	0.845	0.702	0.512	0.590	61.618	7.680
	GPT-4o-mini	0.433	0.479	0.947	0.548	0.668	0.543	0.445	81.398	7.819
	Llama-3.1-70B	0.440	0.579	0.922	0.888	0.662	0.455	0.470	52.570	6.542
	Llama-3.1-8B	0.443	0.586	0.938	0.859	0.673	0.441	0.465	49.986	7.02
MultiUN	Human	n/a	0.312	0.942	0.871	0.612	0.334	0.539	57.357	7.902
	GPT-4o-mini	0.422	0.455	0.961	0.629	0.621	0.518	0.455	74.459	7.875
	Llama-3.1-70B	0.447	0.546	0.912	0.914	0.622	0.329	0.446	52.920	6.903
	Llama-3.1-8B	0.446	0.557	0.950	0.900	0.606	0.295	0.439	55.186	7.014
Wikipedia	Human	n/a	0.286	0.969	0.812	0.723	0.286	0.547	78.316	7.747
	GPT-4o-mini	0.388	0.428	0.963	0.640	0.690	0.446	0.462	78.149	7.696
	Llama-3.1-70B	0.445	0.557	0.931	0.830	0.670	0.278	0.430	57.310	6.681
	Llama-3.1-8B	0.444	0.489	0.963	0.822	0.691	0.322	0.462	58.038	6.246
GovReport	Human	n/a	0.296	0.956	0.815	0.670	0.361	0.548	52.627	6.720
	GPT-4o-mini	0.402	0.420	0.972	0.639	0.683	0.529	0.480	81.374	7.72
	Llama-3.1-70B	0.454	0.511	0.923	0.874	0.667	0.406	0.472	49.294	6.57
	Llama-3.1-8B	0.459	0.498	0.972	0.871	0.650	0.415	0.482	52.712	6.987
Summscreen	Human	n/a	0.308	0.930	0.732	0.557	0.414	0.514	68.971	7.733
	GPT-4o-mini	0.314	0.364	0.974	0.778	0.511	0.345	0.449	61.149	6.667
	Llama-3.1-70B	0.417	0.436	0.984	0.849	0.505	0.281	0.437	49.294	6.413
	Llama-3.1-8B	0.402	0.501	0.987	0.855	0.506	0.295	0.428	62.788	6.067
Seahorse	Human	n/a	0.271	0.952	0.811	0.651	0.518	0.576	60.999	7.067
	GPT-4o-mini	0.353	0.353	0.977	0.786	0.632	0.542	0.517	72.321	7.720
	Llama-3.1-70B	0.417	0.491	0.963	0.831	0.660	0.480	0.489	60.241	6.440
	Llama-3.1-8B	0.402	0.474	0.968	0.839	0.642	0.476	0.490	60.415	7.080

Table 8: Performance on the 10-percent summarization task per dataset

DATASET	Model	R-L(\uparrow)	REP-3(\downarrow)	CoLA \uparrow	COH-2 \uparrow	SH-4 \uparrow	SH-5 \uparrow	AVG \uparrow	HE \uparrow	HUM \uparrow
LCFO.20%										
LexGLUE	Human	n/a	0.455	0.940	0.842	0.688	0.417	0.559	65.036	7.800
	GPT-4o-mini	0.516		0.5822	0.9446	0.5874	0.6286	0.4820	0.4121	8.093
	Llama-3.1-70B	0.507	0.702	0.927	0.860	0.632	0.369	0.417	54.710	7.027
	Llama-3.1-8B	0.501	0.824	0.943	0.882	0.637	0.410	0.409	49.870	6.973
BookSum	Human	n/a	0.344	0.918	0.766	0.621	0.349	0.517	89.376	7.605
	GPT-4o-mini	0.355	0.385	0.975	0.831	0.573	0.441	0.487	95.169	7.432
	Llama-3.1-70B	0.455	0.544	0.956	0.842	0.511	0.314	0.416	69.068	6.296
	Llama-3.1-8B	0.453	0.634	0.971	0.842	0.550	0.394	0.425	76.472	6.457
SQuALITY	Human	n/a	0.395	0.919	0.782	0.565	0.339	0.505	61.257	5.800
	GPT-4o-mini	0.382	0.425	0.969	0.774	0.518	0.328	0.433	79.698	6.720
	Llama-3.1-70B	0.412	0.498	0.963	0.797	0.454	0.205	0.384	41.584	4.027
	Llama-3.1-8B	0.426	0.601	0.979	0.835	0.469	0.233	0.383	47.011	5.587
FacetSum	Human	n/a	0.415	0.940	0.829	0.745	0.490	0.584	72.317	8.147
	GPT-4o-mini	0.477	0.483	0.953	0.685	0.659	0.526	0.468	80.937	7.711
	Llama-3.1-70B	0.474	0.622	0.944	0.899	0.705	0.507	0.487	55.197	6.501
	Llama-3.1-8B	0.456	0.565	0.952	0.894	0.698	0.471	0.490	46.081	6.813
JRC-Acquis	Human	n/a	0.435	0.971	0.860	0.713	0.566	0.605	72.317	7.902
	GPT-4o-mini	0.513	0.566	0.952	0.551	0.685	0.578	0.440	75.351	7.681
	Llama-3.1-70B	0.493	0.792	0.854	0.884	0.648	0.422	0.403	38.876	6.653
	Llama-3.1-8B	0.490	0.788	0.929	0.882	0.633	0.446	0.420	49.870	6.833
MultiUN	Human	n/a	0.422	0.942	0.875	0.605	0.317	0.531	64.540	8.139
	GPT-4o-mini	0.484	0.604	0.954	0.623	0.615	0.482	0.414	72.007	7.917
	Llama-3.1-70B	0.483	0.654	0.907	0.918	0.625	0.289	0.417	37.311	6.694
	Llama-3.1-8B	0.512	0.665	0.925	0.908	0.603	0.326	0.419	45.155	7.028
Wikipedia	Human	n/a	0.355	0.967	0.817	0.738	0.333	0.557	81.923	7.653
	GPT-4o-mini	0.471	0.516	0.960	0.596	0.687	0.432	0.432	76.772	7.609
	Llama-3.1-70B	0.467	0.779	0.877	0.849	0.638	0.300	0.377	55.349	6.493
	Llama-3.1-8B	0.476	0.701	0.940	0.786	0.606	0.257	0.378	51.110	6.014
GovReport	Human	n/a	0.397	0.954	0.823	0.712	0.405	0.563	60.368	8.027
	GPT-4o-mini	0.488	0.521	0.968	0.605	0.684	0.521	0.451	76.887	7.680
	Llama-3.1-70B	0.489	0.638	0.916	0.872	0.634	0.425	0.442	43.421	6.347
	Llama-3.1-8B	0.479	0.531	0.971	0.881	0.623	0.384	0.466	40.544	7.093
Summscreen	Human	n/a	0.395	0.939	0.739	0.552	0.414	0.513	63.691	8.373
	GPT-4o-mini	0.347	0.443	0.969	0.756	0.503	0.346	0.426	60.306	6.200
	Llama-3.1-70B	0.432	0.487	0.979	0.845	0.502	0.294	0.426	60.564	6.627
	Llama-3.1-8B	0.432	0.519	0.987	0.851	0.522	0.322	0.433	49.467	6.413
Seahorse	Human	n/a	0.336	0.964	0.828	0.673	0.533	0.586	66.028	7.907
	GPT-4o-mini	0.415	0.439	0.970	0.722	0.607	0.507	0.474	71.568	8.173
	Llama-3.1-70B	0.454	0.589	0.957	0.834	0.615	0.441	0.452	54.075	6.600
	Llama-3.1-8B	0.469	0.642	0.960	0.847	0.601	0.456	0.444	54.142	6.760

Table 9: Performance on the 20-percent summarization task per dataset.