

Benchmarking LLMs on the Semantic Overlap Summarization Task

Anonymous ACL submission

Abstract

Semantic Overlap Summarization (SOS) is a constrained multi-document summarization task, where the constraint is to capture the common/overlapping information between two alternative narratives. While recent advancements in Large Language Models (LLMs) have achieved exceptional performance in numerous summarization tasks, a benchmarking study of the SOS task using LLMs is yet to be performed. As LLMs’ responses are highly sensitive to variations in prompt design, a major challenge in conducting such a benchmarking study is to systematically explore a variety of prompts before drawing a reliable conclusion. Fortunately, the TELeR taxonomy has been recently proposed, which can be used to design and explore various prompts for LLMs. Using this TELeR taxonomy, this paper comprehensively evaluates 16 popular LLMs on the SOS Task. We evaluate and report on 905,216 LLM generated summaries using well-established metrics like ROUGE, BERTscore, and SEM- F_1 on two different datasets of alternative narratives and we also conduct human evaluation on 540 of those summaries for further analysis. We conclude the paper by analyzing the strengths and limitations of various LLMs in terms of their capabilities in capturing overlapping information¹.

1 Introduction

Large Language Models (LLMs) represent a groundbreaking advancement in the research landscape of Natural Language Processing (NLP) and Artificial Intelligence (AI). Trained on large bodies of text data, LLMs excel in generating coherent and human-like text. These models have been evaluated in a wide range of NLP tasks (Bubeck et al., 2023; Dai et al., 2022; Du et al., 2022; Smith et al.,

¹The code and datasets used to conduct this study are available at https://anonymous.4open.science/r/llm_eval-E16D

2022) across several areas, including software development, law, and medicine (Schäfer et al., 2024; School, 2023; Thirunavukarasu et al., 2023). However, there are still areas and tasks where LLMs are yet to be rigorously evaluated. One such task is Semantic Overlap Summarization (SOS) (Bansal et al., 2022c; Karmaker Santu et al., 2018), where the goal is to summarize the common/overlapping information between two alternative narratives.

In this paper, we conduct a comprehensive benchmarking study of the SOS task using 16 popular LLMs. Conducting such a benchmarking study is challenging because of the large variations in LLMs’ performance when different prompt types/styles are used and different degrees of detail are provided in the prompts. Indeed, Rodriguez et al. (2023) shows varying performance on the CM1 dataset (Hayes et al., 2006) across many different prompts with F1-scores ranging from 0.21 to 0.54, exhibiting a 0.33 point difference. To address this issue, Santu and Feng (2023) recently proposed a general taxonomy that can be used to design diverse prompts with specific properties in order to perform a wide range of complex tasks. Using this TELeR taxonomy, we devised a comprehensive set of prompts with different degrees of detail to perform the SOS task on two different alternative narratives datasets. One dataset is the previously introduced AllSides dataset released by Bansal et al. (2022c), and the second one is our original contribution with extensive human annotation efforts, which we name the *PrivacyPolicyPairs* (3P) dataset.

Figure 1 illustrates an example of the SOS task when comparing two alternative privacy policies in the 3P dataset, where the green text denotes the output (common information) from two input privacy policies, one from Google and one from Apple. In this case, we have two competing platforms that provide similar types of services (e.g. Cloud Storage or Streaming Services) and each company



Figure 1: Example of a SOS task using two alternative privacy policy narratives.

lays out its practices when handling your private information. These documents contain important information regarding your privacy but can be long and cumbersome to read. The SOS task can help users by briefly identifying the common practices followed by each company.

For evaluation, we report well-established metrics like ROUGE, BERTscore, and SEM- F_1 on Allsides and 3P datasets for each combination of LLMs and prompt style, totaling 905,216 unique samples for analysis. We further evaluate a subset of samples using human annotators to truly gauge the capabilities of LLMs in capturing and synthesizing overlapping information from multiple narratives. We conclude the paper by analyzing the strengths and limitations of various LLMs for the same task.

2 Related Work

Text Summarization: SOS is essentially a summarization task. Over the past two decades, many document summarization approaches have been investigated (Zhong et al., 2019). The two most popular among them are *extractive* approaches (Cao et al., 2018; Narayan et al., 2018; Wu and Hu, 2018; Zhong et al., 2020) and *abstractive* approaches (Bae et al., 2019; Liu et al., 2017; Nallapati et al., 2016). Some researchers have tried combining extractive and abstractive approaches (Chen and Bansal, 2018; Hsu et al., 2018; Zhang et al., 2019).

The SOS Task: *Semantic Overlap Summarization* can be framed as a multi-document summarization task, *i.e.*, multi-seq-to-seq task (Goldstein et al., 2000; Yasunaga et al., 2017; Zhao et al., 2020; Ma et al., 2020; Meena et al., 2014; Lebanoff et al., 2018; Fabbri et al., 2019). However, unlike typical multi-document summarizing tasks, *SOS* aims to summarize multiple alternative narratives with

an overlapping constraint (Bansal et al., 2022c), *i.e.*, the output should only contain the common information from both input narratives (Santu et al., 2018). The availability of data for this task is relatively small so recently Bansal et al. (2022a) proposed a method for training models by utilizing synthetically generated data.

Transformers and LLMs: Encoder-decoder-based transformer models have recently gained a lot of attraction, especially for abstractive summarization tasks, (Rush et al., 2015; Chopra et al., 2016; Zhou et al., 2017; Paulus et al., 2017). Training a generic language model on a large corpus of data and then transferring/fine-tuning it for the summarization job has become a standard approach (Radford et al., 2019; Devlin et al., 2019; Lewis et al., 2019; Xiao et al., 2020; Yan et al., 2020; Zhang et al., 2019; Raffel et al., 2019). Transformer architecture features more parallelizable training, better-scaling properties, and a built-in attention mechanism, allowing large language models (LLMs) to emerge. Made up of billions of parameters, many LLMs like GPT-2 (Radford et al.), and LLaMA (Touvron et al., 2023a) have showcased strong abilities at generating text. Then, with the introduction of Reinforcement Learning From Human Feedback (Ouyang et al., 2022), LLMs became even more powerful, allowing users to interact with them as data with natural language queries.

Prompt Engineering for LLMs: “Prompt Engineering” is a technique for maximizing the utility of LLMs in various tasks (Zhou et al., 2022). It involves crafting and revising the query or context to elicit the desired response or behavior from LLMs (Brown et al., 2022). Prompt engineering is an iterative process requiring multiple trial and error runs (Shao et al., 2023). In fact, differences in prompts along several key factors can significantly

155 impact the accuracy and performance of LLMs in
156 complex tasks. To address this issue, [Santu and](#)
157 [Feng \(2023\)](#) recently proposed the TELeR taxon-
158 omy, which can serve as a unified standard for
159 benchmarking LLMs’ performances by exploring
160 a wide variety of prompts in a structured manner.

161 3 Background

162 3.1 Semantic Overlap Summarization Task

163 Semantic Overlap Summarization (SOS) is a task
164 aimed at extracting and condensing shared infor-
165 mation between two input documents, D_A and D_B .
166 The output, denoted as D_O , is generated in natu-
167 ral language and only includes information present
168 in both input documents. The task is framed as a
169 constrained multi-seq-to-seq (text generation) task,
170 where brevity is emphasized to minimize the rep-
171 etition of overlapping content. The output can be
172 extractive summaries, abstractive summaries, or a
173 combination of both ([Karmaker Santu et al., 2018](#)).

174 Furthermore, SOS adheres to the commutative
175 property, meaning the order of input documents
176 doesn’t affect the output summary; $D_A \cap_O D_B =$
177 $D_B \cap_O D_A$. To facilitate research in this area,
178 [Bansal et al. \(2022c\)](#) introduced the AllSides
179 dataset for training and evaluation, which we also
180 used for evaluation in this work.

181 3.2 Prompting With TELeR Taxonomy

182 As shown in Figure 6, the TELeR taxonomy in-
183 troduced by [Santu and Feng \(2023\)](#) categorizes
184 complex task prompts based on four criteria.

- 185 1. **Turn:** This refers to the number of turns or shots
186 used while prompting an LLM to accomplish
187 a complex task. In general, prompts can be
188 classified as either single or multi-turn.
- 189 2. **Expression:** This refers to the style of expres-
190 sion for interacting with the LLM, such as ques-
191 tioning or instructing.
- 192 3. **Level of Details:** This dimension of prompt
193 style deals with the granularity or depth of ques-
194 tion or instruction. Prompts with higher levels
195 of detail provide more granular instructions.
- 196 4. **Role:** LLMs can provide users with the option
197 of specifying the role of the system. The re-
198 sponse of LLM can vary due to changes in role
199 definitions in spite of the fact that the prompt
200 content remains unchanged.

201 The taxonomy outlines 7 distinct levels starting
202 from level 0 to level 6. With each increase in level
203 comes an increase in complexity of the prompt.

In level 0, only data/context is provided with no
further instruction. Level 1 extends level 0 by pro-
viding single-sentence instruction. Then level 2
extends level 1, and so on, until level 6, where
all characteristics of previous levels are provided
along with the additional instruction for the LLM to
explain its output. For more details on the TELeR
taxonomy and its applications, see [Santu and Feng](#)
[\(2023\)](#). For convenience, we include the outline
diagram from the paper in Appendix A.2.

214 4 The Benchmark Datasets

215 4.1 The AllSides Data

216 The AllSides dataset is collected from All-
217 Sides.com, a third-party online news forum known
218 for presenting news and information from vari-
219 ous political perspectives. [Bansal et al. \(2022c\)](#)
220 crawled news articles from AllSides.com to build
221 the dataset, focusing on narratives covering 2, 925
222 events. These articles provide contrasting perspec-
223 tives from media outlets affiliated with “Left” and
224 “Right” political leanings. Additionally, each event
225 includes a factual description labeled as “Theme”
226 by AllSides, which can serve as a neutral perspec-
227 tive for readers (ground truth for common facts).

228 The test set is comprised of 137 narrative pairs
229 where each sample has 4 references: one from
230 AllSides and 3 from human annotators, totaling
231 548 reference summaries for the test set.

232 4.2 The *PrivacyPolicyPairs* (3P) Data

233 For a more comprehensive evaluation, we introduce
234 the *PrivacyPolicyPairs* (3P) dataset, an additional
235 evaluation set for the SOS task, which contains
236 135 human annotated samples. Our (3P) dataset is
237 built on the OPP-115 Corpus introduced by [Wilson](#)
238 [et al. \(2016\)](#), which comprises 115 privacy policies
239 (267K words) spanning 15 sectors (Arts, Shopping,
240 News, etc.). The policy data of the OPP-115 corpus
241 are also tagged with the following categories:

- 242 • First Party Collection/Use
- 243 • Third Party Sharing/Collection
- 244 • User Choice/Control
- 245 • User Access, Edit, & Deletion
- 246 • Data Retention
- 247 • Data Security
- 248 • Policy Change
- 249 • Do Not Track
- 250 • International & Specific Audiences

3P Data Sample		
Category: Data Security		
Policy 1: Amazon (410 Words)	Policy 2: Lids (312 Words)	
<p>Amazon.com knows that you care how information about you is used and shared, and we appreciate your trust that we will do so carefully and sensibly</p> <p>...</p> <p>We work to protect the security of your information during transmission by using Secure Sockets Layer (SSL) software, which encrypts information you input. We reveal only the last four digits of your credit card numbers when confirming an order. Of course, we transmit the entire credit card number to the appropriate credit card company during order processing. It is important for you to protect against unauthorized access to your password and to your computer. Be sure to sign off when finished using a shared computer. Click here for more information on how to sign off</p> <p>...</p>	<p>Any personal information that we collect will be stored in secure servers hosted in the U.S. or Canada</p> <p>...</p> <p>We work to protect the security of your information during transmission by using Thawte Certified Secure Sockets Layer (SSL) software, which encrypts information you input. We reveal only the last four digits of your credit card numbers when confirming an order. Of course, we transmit the entire credit card number to the appropriate credit card company during order processing.</p> <p>Security lies in your hands as well. It is important for you to protect against unauthorized access to your password and to your computer. Be sure to sign off when finished using a shared computer. In the event of unauthorized use of your credit card, you must notify your credit card provider in accordance with its reporting rules and procedures.</p> <p>...</p>	
Reference Summaries		
A_1	A_2	A_3
<p>We work to protect the security of your information during transmission by using Secure Sockets Layer (SSL) software, which encrypts information you input. We reveal only the last four digits of your credit card numbers when confirming an order. Of course, we transmit the entire credit card number to the appropriate credit card company during order processing. It is important for you to protect against unauthorized access to your password and to your computer. Be sure to sign off when finished using a shared computer.</p>	<p>Companies work to protect the security of your information during transmission by using Secure Sockets Layer (SSL) software, which encrypts information you input. They reveal only the last four digits of your credit card numbers when confirming an order. Of course, They transmit the entire credit card number to the appropriate credit card company during order processing. It is important for you to protect against unauthorized access to your password and to your computer. Hence, be sure to sign off when finished using a shared computer.</p>	<p>Even though the entire credit card number is transmitted, only the last 4 digits of the credit card number is visible during confirmation. SSL is used to save info during transmission. Sign off is recommended.</p>

Table 1: A single sample from the 3P dataset. For each sample, you are given the category name, company names, the corresponding policy subsections, the count of words in each policy, and the 3 reference summaries. The highlighted text shows the overlapping information.

• Other

These annotations were also associated with a *text span* in the privacy policy to denote where the labels were relevant.

Our motivation behind introducing a new dataset for SOS evaluation is the following: 1) it extends the amount of available testing data from just 137 samples from the AllSides evaluation set to 272 total evaluation samples with a combined total of 953 human annotations for the two datasets; 2) The 3P dataset represents a new type of documents in the form of semi-structured privacy policies as opposed to the news articles that make up the AllSides data; 3) News datasets are abundant, and LLMs are extensively pretrained on them in comparison to relatively infrequent privacy policy data; hence the 3P dataset is supposedly more challenging for LLMs.

Constructing the 3P Dataset: The 3P dataset includes pairs of passages taken from the OPP-115 corpus and tasks the annotator with finding the semantically overlapping information between them. A data sample is shown in Table 1. Each sample comprises 2 source documents (two alternative privacy policy narratives), the category they fall under,

and 3 reference overlap summaries. The company names and word counts are also included.

When curating this dataset, we wanted to ensure each passage pair had some degree of overlap. To facilitate this goal, we reversed the process followed by the original authors and grouped the documents back into their respective sectors. Then, we built pairs of passages for each document in each sector according to the categories they were originally labeled with. This process resulted in 6110 passage pairs across all sectors.

3P Dataset Statistics	
# Samples	135
Avg. # Words per Document	331.00
Avg. # Words per Document Pair	662.01
Avg. # Sentences per Document	14.96
Avg. # Sentences per Document Pair	28.99
Avg. # Words per Reference	22.46
Avg. # Sentences per Reference	1.75

Table 2: Dataset statistics for the 3P dataset consisting of 135 document pairs with 3 references each.

Of the sectors, we chose to focus on three: *eCommerce*, *Technology*, and *Food and Drink*, due to their popularity as well as diversity among each other. From these sectors, we collected 346 pas-

sage pairs to annotate. For this task, three human annotators were asked to write a summary of common information present in each document pair. Conflicting summaries arose when there was no overlap or the annotators considered shared words as overlap. To address these issues, we retain only the policy pairs where at least two annotators wrote at least 15 words as their reference summaries. After annotating, resolving conflicts, and removing samples with no overlap, the process yielded us 3 annotations per passage pair for a total of 405 annotations for 135 high-quality samples. The final dataset statistics are listed in Table 2.

5 Methodology

5.1 Large Language Models Evaluated

We choose to test our datasets using 7 families of instruction-tuned LLMs, totaling 16 models. All evaluated models are listed in Table 3. For the commercial LLMs (OpenAI and Google), we used their provided APIs for summary generation, but for open-source LLMs, we used the huggingface transformers library (Wolf et al., 2020) to access model weights and perform generation on a server with 4X A4500 20GB GPUs. For additional speedup, we leveraged the vLLM library (Kwon et al., 2023).

LLM Family	Model
Google Gemini (Team et al., 2024)	gemini-1.5-pro-001 (May 2024)
OpenAI (OpenAI, 2023)	gpt-3.5-turbo-0125 (May 2024)
MosaicML MPT (Team, 2023)	mosaicml/mpt-7b-chat (7B) mosaicml/mpt-30b-chat (30B) mosaicml/mpt-7b-instruct (7B) mosaicml/mpt-30b-instruct (30B)
LMSYS Vicuna (Zheng et al., 2023)	lmsys/vicuna-7b-v1.5 (7B) lmsys/vicuna-13b-v1.5 (13B) lmsys/vicuna-7b-v1.5-16k (7B) lmsys/vicuna-13b-v1.5-16k (13B)
MistralAI (Jiang et al., 2023)	mistralai/Mistral-7B-Instruct-v0.1 (7B) mistralai/Mistral-7B-Instruct-v0.2 (7B)
MetaAI Llama2 (Touvron et al., 2023b)	meta-llama/Llama-2-7b-chat-hf (7B) meta-llama/Llama-2-13b-chat-hf (13B)
Microsoft Phi-3 (Abdin et al., 2024)	microsoft/Phi-3-mini-4k-instruct (3.8B) microsoft/Phi-3-mini-128k-instruct 3.8B)

Table 3: The list of models evaluated in this paper. We use 7 families of models, 2 of which are closed source, and 5 open source. Parameter counts of open source models are included in parentheses. OpenAI and Google have not reported the parameter counts of their models.

We prompted LLMs in a zero-shot setting with TELeR as zero-shot approaches to NLP tasks have gained popularity with the growing capabilities of LLMs. For example, works from Sarkar et al. (2023, 2022) explore their zero-shot use cases in topic inference and text classification. For this

study, we used TELeR levels 0 through 4 (5 out of the 7). We chose not to prompt using levels 5 and 6 because their use of retrieval augmented prompting does not necessarily apply to the SOS task due to all relevant context being present, *i.e.*, the two source narratives are already provided as part of the prompt. Furthermore, requirement number 5 for level 6 also specifies asking the LLM to explain its own output, which would negatively affect the generated summaries during evaluation. We also experiment with in-context learning prompts (Brown et al., 2020).

5.2 Designed Prompts

For each template, we use the following outline for our prompt design.

- **TELeR Level 0: {Document 1} {Document 2}**
- **TELeR Level 1:**
 - Document 1: **{Document 1 Content}**
 - Document 2: **{Document 2 Content}**
 - Summarize the overlapping information between these two documents
- **TELeR Level 2:**
 - {TELeR Level 1 Prompt Text}**
 - This information must keep in mind the 5W1H facets of the documents. Do not include any uncommon information.
- **TELeR Level 3:**
 - {TELeR Level 1 Prompt Text}**
 - This information must keep in mind the 5W1H facets of the documents.
 - Do not include uncommon information.
- **TELeR Level 4:**
 - {Level 3 Prompt Text}**.
 - Your response will be evaluated against a set of reference summaries. Your score will depend on how semantically similar your response is to the reference.
- **In-context Learning:**
 - Document 1: **{Example Doc. 1 Content}**
 - Document 2: **{Example Doc. 2 Content}**
 - Summary: **{Example Summary}**

 - Document 1: **{Document1 Content}**
 - Document 2: **{Document2 Content}**
 - Summary:

To ensure comprehensive prompt engineering, we created groups of templates for TELeR levels 0 through 4, In-Context Learning (Brown et al., 2020) formats, and also for system roles. In each template group, we create variations of prompts that follow their respective formats. For example, the group of TELeR L1 prompts is comprised of 5 general prompts, 3 AllSides-specific prompts, and 3 3P-specific prompts. Then, to construct our final set of prompts, we took all possible combinations of system roles and prompts. A breakdown of the

variation counts for each group is shown in Table 4. Using this prompting strategy, we’ve created 56,576 unique prompts for each of our 16 evaluated LLMs, totaling 905,216 evaluation samples. See appendix A.2 for the exact prompts that were used.

Template Group	For PPP	For AllSides	For Both	Total
System Role	2	2	6	10
TELeR L0	0	0	1	1
TELeR L1	3	3	5	11
TELeR L2	3	3	3	9
TELeR L3	3	3	2	8
TELeR L4	3	3	2	8
In-Context Learning	0	0	1	1

Table 4: The number of prompts created for each template group. The "For PPP/AllSides" columns indicate how many prompts were created for that dataset only. The "For Both" column is for the prompts that could be applied to both datasets. For exact prompt details, refer to Appendix A.2 for exact prompt contents.

5.3 Evaluation

5.3.1 Automatic Evaluation Metrics

ROUGE: ROUGE (Lin, 2004) is a family of metrics that score the lexical overlap between the generated text and the reference text. We used 3 variations, R-1, R-2, and R-L, which are widely adopted for evaluating text summarizing tasks. However, despite its popularity, works like Akter et al. (2022b) and Bansal et al. (2022b) show that ROUGE is an unsuitable metric for comparing semantics.

BERTscore: BERTscore is a metric that utilizes contextual embeddings from transformer models like BERT to evaluate the semantic similarity between the generated text and reference text. For this study, we compute BERTscore with the hashcode `roberta-large_L17_no-idf_version=0.3.12(hug_trans=4.40.2)-rescaled`.

SEM-F1: While ROUGE and BERTscore are useful and powerful metrics, SEM-F1 was specifically designed for the SOS task. SEM-F1 leverages rigorously fine-tuned sentence encoders to evaluate the SOS task using sentence-level similarity unlike BERTscore which utilizes token-level similarity. For this study, we compute SEM-F1 with 3 underlying models: USE (Cer et al., 2018), RoBERTa (Zhuang et al., 2021), and DistilRoBERTa (Sanh et al., 2019).

5.3.2 Human Evaluation

For our human evaluation strategy, we recruited 3 volunteer annotators. These annotators evaluated 15 randomly chosen dataset samples, where 7

were picked from the AllSides dataset and 8 were picked from the 3P dataset. For each dataset sample, if we consider all prompts for every possible template combination and every model, it will amount to 3,328 annotations required for each of the 15 dataset samples, making human judgment very time-consuming. To solve this challenge, we reduce this number by 1) choosing a subset of models to evaluate and 2) choosing the best-evaluated prompts for each of the 6 template groups based on their average performance in terms of automated metrics on both datasets. This strategy reduced the number of summaries per sample from 3,328 to 36 summaries per sample, giving us a total of 540 annotations per human annotator. The humans were tasked to score the summaries on a scale of 0-5 based on how well they captured the overlapping information of the 2 given source documents. After individually scoring the summaries, the annotators sat together to negotiate a final score to assign to each sample, giving us 2160 annotation scores across all samples.

6 Results

For clarity, we show our results for automatic evaluation scores on the largest or newest models of each family in Figure 2. This Figure shows the highest scores achieved by each model over the set of given TELeR prompts. The first two columns are for the commercial LLMs (GPT-3.5-Turbo and Gemini-Pro), and the 5 columns on the right show our open-source LLMs. In general, we observe that GPT-3.5-Turbo and Gemini-Pro beat the open-source models across all benchmarks but we also find that Mistral-7B-Instruct-v0.2 exhibits competitive performance across all metrics despite being only a 7-billion parameter model. For a comprehensive breakdown of the model scores achieved by each LLM and for each metric, refer to appendix (Figure 4).

Annotator/Metric Agreement In Figure 3, we show Pearson correlation and Kendall’s τ correlation to compare human annotator scores with our automatic evaluation metrics. We denote the final annotation score by Ann_{comb} . From this table, all metrics have a relatively low correlation with human judgments, again demonstrating the limitations of automated metrics for evaluating text generation. Interestingly, we see that SEM-F1 best correlates with the human annotators, demonstrating its superior quality over ROUGE and BERTscore

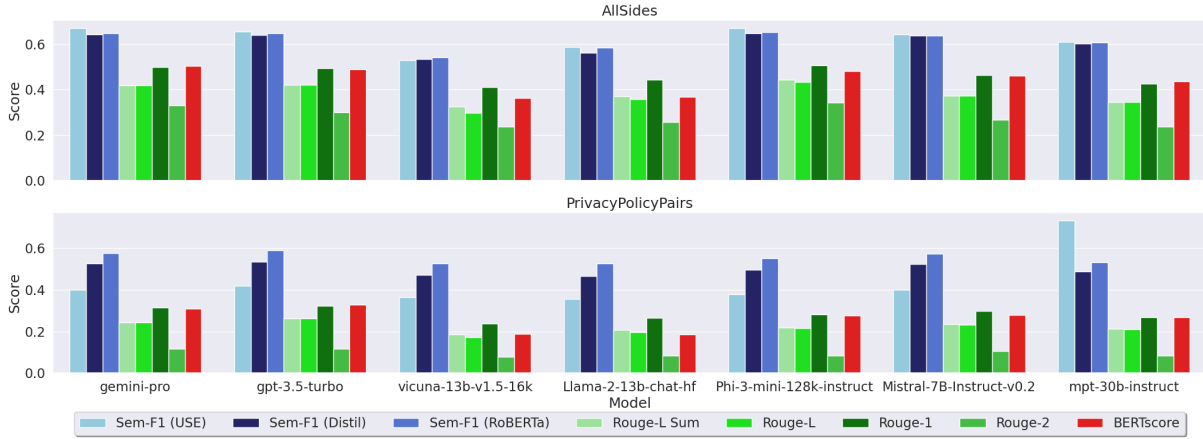


Figure 2: Best scores over each TELER prompt level for the largest model of each family of LLMs and for each dataset. Red shows BERTscore, green shows ROUGE, and blue shows Sem-F1. A full breakdown of max scores obtained by each model is shown in Appendix A.1

Pearson Correlation - Annotators vs. Metrics

Ann ₁	0.2	0.29	0.31	0.09	0.08	0.1	0.05	0.24
Ann ₂	0.26	0.35	0.42	0.19	0.19	0.19	0.1	0.34
Ann ₃	0.06	0.1	0.11	0	0	0.01	-0.02	0.06
Ann _{comb}	0.15	0.24	0.25	0.06	0.05	0.06	0.02	0.2

Kendall's τ - Annotators vs. Metrics

Ann ₁	0.05	0.07	0.05	0.01	0.01	0	0	0.05
Ann ₂	0.09	0.11	0.1	0.09	0.09	0.09	0.06	0.09
Ann ₃	0	0.04	0.04	-0.04	-0.04	-0.05	-0.03	-0.01
Ann _{comb}	0.03	0.05	0.03	0.01	0	-0.01	-0.01	0.03

SEM-F1_{USE} SEM-F1_{Distil} SEM-F1_{Rob} ROUGE-L_{SUM} ROUGE-L ROUGE-1 ROUGE-2 BERTscore

Figure 3: Pearson correlation and Kendall's τ scores between annotator scores and automatic evaluation metrics (higher is better). The "comb" subscript shows the combined score where the annotators sat with each other to settle on a final score for each annotation sample.

for multi-document summary evaluation.

AllSides Vs. 3P: In Table 5, we show average scores across all LLMs for each template prompt. The highest scores for each column are bolded. From this table, it is evident that scores on the 3P dataset tend to trail significantly behind the AllSides dataset. Additionally, we observe that for automatic evaluation, TELER level 1 consistently performs the best across the AllSides dataset, while on the 3P dataset, the SEM-F1 scores suggest better performance in TELER level 2 and 4.

Human Preference on Model and Template:

While Table 5 shows that the automatic evaluations tend to have a preference towards TELER L1 prompts, Table 6 shows that human annotators actually tend to prefer TELER L2 prompts instead. However, this preference is only 0.04 points ahead of the next best. The table also indicates the annotators' preference towards gpt-3.5-turbo for the commercial LLMs. Then, for the open-source LLMs, mpt-30b-chat

was the most preferred, with an average annotator score of 3.39. However, it is important to note that Phi-3-mini-128k-instruct and Mistral-7B-Instruct-v0.2 match and beat gemini-pro, respectively, according to humans.

7 Discussion

The main takeaways from our study are as follows:

Finding-1: "3P" dataset is harder than "AllSides" for LLMs in the context of SOS task.

To elaborate, Table 5 shows a clear difference in scores between the AllSides data and the 3P data. These differences can possibly be explained for the following reasons. The average document word count for both datasets has a significant difference but is well within the context windows of LLMs. For the AllSides data, the average is 504.51 while for the 3P data, it's 662.01. Another difference worth noting is the amount of overlapping

Dataset	Template	BERTscore	R-1	R-2	R-L Sum	R-L	Sem-F1 (Distil)	Sem-F1 (RoBERTa)	Sem-F1 (USE)
AllSides	ICL	0.453	0.46	0.267	0.367	0.367	0.621	0.639	0.651
	L0	0.391	0.399	0.209	0.315	0.291	0.6	0.618	0.614
	L1	0.503	0.507	0.342	0.442	0.433	0.646	0.652	0.671
	L2	0.437	0.466	0.278	0.388	0.369	0.631	0.636	0.653
	L3	0.461	0.465	0.273	0.377	0.376	0.639	0.646	0.658
Privacy Policy Pairs (3P)	L4	0.477	0.467	0.268	0.376	0.376	0.641	0.647	0.656
	ICL	0.262	0.278	0.092	0.219	0.219	0.499	0.55	0.371
	L0	0.138	0.226	0.067	0.187	0.164	0.509	0.567	0.399
	L1	0.329	0.324	0.118	0.262	0.262	0.535	0.588	0.419
	L2	0.267	0.307	0.109	0.254	0.234	0.531	0.589	0.414
L3	0.256	0.278	0.08	0.211	0.21	0.517	0.578	0.385	
	L4	0.299	0.314	0.112	0.244	0.243	0.535	0.577	0.734

Table 5: Average scores per metric broken down by level and dataset. TELeR Levels are denoted by "Lx" and In-Context Learning is denoted by "ICL". The highest of each metric and dataset are in bold.

Model	Score (0-5)
gemi-ni-pro	3.37
gpt-3.5-turbo	3.53
mpt-30b-chat	3.39
Mistral-7B-Instruct-v0.2	3.38
Phi-3-mini-128k-instruct	3.37
vicuna-13b-v1.5-16k	3.32
Template	Score (0-5)
ICL	3.08
TELeR L1	3.38
TELeR L2	3.42
TELeR L3	3.32
TELeR L4	3.32

Table 6: Average negotiated preference score for each model and prompt template. "ICL" represents the In-Context Learning style prompts while "Lx" refers to the level of TELeR prompt.

tokens present in each dataset. Utilizing the NLTK library (Bird et al., 2009) for tokenization, we find that the AllSides dataset has an average of 23.36% unique token overlap between source documents while it's 17.95% only for 3P. Aside from the numerical differences, it is also important to note the compositional differences between these datasets. The AllSides data is comprised of news articles that cover the same events. By the nature of covering the same events, it is more likely for overlapping spans of text to occur in the form of quotes or headlines. In contrast, although the documents of the 3P dataset are similar in terms of their document structure and the types of excerpts they fall under, but they are not essentially written about the same subjects. Each document pair in the 3P dataset represents two different companies whose policies can differ greatly while using similar language.

Finding-2: SEM-F1 remains the recommended evaluation metric for the SOS task.

SEM-F1 showcases the best correlation with human annotators, while ROUGE continues to show its limitations, which is consistent with the previous

findings from the literature (Akter et al., 2022a). Figure 3 supports these claims, showing correlations between annotators and ROUGE that even go into the negatives. This figure also further highlights the work that needs to be done in order to improve automatic evaluation to the point where we can rely on them more seriously and let go of expensive human evaluation.

Finding-3: Use gpt-3.5-turbo with TELeR L2 prompts for best results in the SOS task

As of the time of writing, closed-source commercial LLMs remain the top performers in text generation. However the quality of open-source models are not too far off from gpt-3.5-turbo according to our human evaluators as noted in Table 6 and most of our evaluated models even compete with gemini-pro. Aside from model preferences it is also important to note that In-Context Learning styled prompts have been shown as a less effective prompting method compared to TELeR in the constrained multi-document summarization setting.

8 Conclusion

In this study, we provide a comprehensive look into the capability of LLMs for the Semantic Overlap Summarization (SOS) task. To facilitate robust evaluation, we test on a previously created dataset and additionally introduce the *PrivacyPolicyPairs* (3P) dataset. We use the TELeR prompting taxonomy to devise a set of hand-crafted prompts that generate the highest scores we could achieve with pre-trained instruction-tuned LLMs and found that: 1) the 3P dataset is a harder benchmark for LLMs 2) SEM-F1 is still the best method for evaluating on SOS but far from ideal and 3) based on our testing methodology, the best summarization results in a zero-shot setting can be accomplished using gpt-3.5-turbo with TELeR L2 style prompts.

9 Limitations

The primary limitation of this work is the size of the dataset. At only 135 samples, it is not feasible to train a model on just the 3P data alone. However while the size of the new dataset is small, there is a large amount of time and resource that is required to build a dataset of this nature. Firstly, this dataset requires that for each sample, we find two documents that share an overlapping narrative. Second, each sample is annotated manually by 3 people which for this dataset results in 405 annotations. That is without considering the other annotations where no overlap was found. Third, there have been several instances where disagreements need to be resolved which requires further discussion among annotators. Despite these limitations it is worth noting that this work effectively doubles the amount of samples to evaluate on the SOS task when considering both AllSides data and 3P data combined. In the future, a larger scale effort will be needed to increase the space of data for the SOS task.

Another limitation is that we did not perform any fine-tuning on these models. All scores were obtained using the pre-trained weights for each model. This means that it's possible for additional performance to be gained using methods like LoRA (Hu et al., 2021). However the main goal of this study was to benchmark LLMs to set new baselines for the SOS task. In that regard we believe this to be an appropriate setup.

References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Qin Cai, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Yen-Chun Chen, Yi-Ling Chen, Parul Chopra, Xiyang Dai, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Victor Fragoso, Dan Iter, Mei Gao, Min Gao, Jianfeng Gao, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Ce Liu, Mengchen Liu, Weishung Liu, Eric Lin, Zeqi Lin, Chong Luo, Piyush Madan, Matt Mazzola, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas

Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Xin Wang, Lijuan Wang, Chunyu Wang, Yu Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Haiping Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Sonali Yadav, Fan Yang, Jianwei Yang, Ziyi Yang, Yifan Yang, Donghan Yu, Lu Yuan, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lina Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). (arXiv:2404.14219). ArXiv:2404.14219 [cs].

Mousumi Akter, Naman Bansal, and Shubhra Kanti Karmaker. 2022a. [Revisiting automatic evaluation of extractive summarization task: Can we do better than ROUGE?](#) In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1547–1560, Dublin, Ireland. Association for Computational Linguistics.

Mousumi Akter, Naman Bansal, and Shubhra Kanti Karmaker Santu. 2022b. [Revisiting automatic evaluation of extractive summarization task: Can we do better than rouge?](#) In *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 1547–1560. Association for Computational Linguistics.

Sanghwan Bae, Taek Kim, Jihoon Kim, and Sang-goo Lee. 2019. Summary level training of sentence rewriting for abstractive summarization. *arXiv preprint arXiv:1909.08752*.

Naman Bansal, Mousumi Akter, and Shubhra Kanti Karmaker Santu. 2022a. [Learning to generate overlap summaries through noisy synthetic data](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, page 11765–11777, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Naman Bansal, Mousumi Akter, and Shubhra Kanti Karmaker Santu. 2022b. [Sem-fl: an automatic way for semantic evaluation of multi-narrative overlap summaries at scale](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, page 780–792, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Naman Bansal, Mousumi Akter, and Shubhra Kanti Karmaker Santu. 2022c. [Semantic overlap summarization among multiple alternative narratives: An exploratory study](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, page 6195–6207, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. " O'Reilly Media, Inc."

672	Hannah Brown, Katherine Lee, Fatemehsadat	<i>the North American Chapter of the Association for</i>	729
673	Mireshghallah, Reza Shokri, and Florian Tramèr.	<i>Computational Linguistics: Human Language Tech-</i>	730
674	2022. What does it mean for a language model to	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	731
675	preserve privacy? In <i>2022 ACM Conference on</i>	4171–4186, Minneapolis, Minnesota. Association for	732
676	<i>Fairness, Accountability, and Transparency</i> , pages	Computational Linguistics.	733
677	2280–2292.		
678	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	Nan Du, Yanping Huang, Andrew M Dai, Simon Tong,	734
679	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun,	735
680	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. 2022.	736
681	Askeell, Sandhini Agarwal, Ariel Herbert-Voss,	Glam: Efficient scaling of language models with	737
682	Gretchen Krueger, Tom Henighan, Rewon Child,	mixture-of-experts. In <i>International Conference on</i>	738
683	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	<i>Machine Learning</i> , pages 5547–5569. PMLR.	739
684	Clemens Winter, Christopher Hesse, Mark Chen,	Alexander R Fabbri, Irene Li, Tianwei She, Suyi Li,	740
685	Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin	and Dragomir R Radev. 2019. Multi-news: A	741
686	Chess, Jack Clark, Christopher Berner, Sam Mc-	large-scale multi-document summarization dataset	742
687	Candlish, Alec Radford, Ilya Sutskever, and Dario	and abstractive hierarchical model. <i>arXiv preprint</i>	743
688	Amodei. 2020. Language models are few-shot learn-	<i>arXiv:1906.01749</i> .	744
689	ers . (arXiv:2005.14165). ArXiv:2005.14165 [cs].		
690	Sébastien Bubeck, Varun Chandrasekaran, Ronen El-	Jade Goldstein, Vibhu O Mittal, Jaime G Carbonell, and	745
691	dan, Johannes Gehrke, Eric Horvitz, Ece Kamar,	Mark Kantrowitz. 2000. Multi-document summariza-	746
692	Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lund-	tion by sentence extraction. In <i>NAACL-ANLP 2000</i>	747
693	berg, et al. 2023. Sparks of artificial general intelli-	<i>Workshop: Automatic Summarization</i> .	748
694	gence: Early experiments with gpt-4. <i>arXiv preprint</i>	J.H. Hayes, A. Dekhtyar, and S.K. Sundaram. 2006. Ad-	749
695	<i>arXiv:2303.12712</i> .	vancing candidate link generation for requirements	750
696	Ziqiang Cao, Wenjie Li, Sujian Li, and Furu Wei. 2018.	tracing: the study of methods . <i>IEEE Transactions on</i>	751
697	Retrieve, rerank and rewrite: Soft template based	<i>Software Engineering</i> , 32(1):4–19.	752
698	neural summarization . In <i>Proceedings of the 56th</i>	Wan-Ting Hsu, Chieh-Kai Lin, Ming-Ying Lee, Kerui	753
699	<i>Annual Meeting of the Association for Computational</i>	Min, Jing Tang, and Min Sun. 2018. A unified model	754
700	<i>Linguistics (Volume 1: Long Papers)</i> , pages 152–161,	for extractive and abstractive summarization using	755
701	Melbourne, Australia. Association for Computational	inconsistency loss. <i>arXiv preprint arXiv:1805.06266</i> .	756
702	Linguistics.		
703	Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua,	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	757
704	Nicole Limtiaco, Rhomni St. John, Noah Constant,	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	758
705	Mario Guajardo-Cespedes, Steve Yuan, Chris Tar,	Weizhu Chen. 2021. Lora: Low-rank adaptation of	759
706	Brian Strope, and Ray Kurzweil. 2018. Universal	large language models .	760
707	sentence encoder for english . In <i>Proceedings of the</i>	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	761
708	<i>2018 Conference on Empirical Methods in Natural</i>	sch, Chris Bamford, Devendra Singh Chaplot, Diego	762
709	<i>Language Processing: System Demonstrations</i> , page	de las Casas, Florian Bressand, Gianna Lengyel,	763
710	169–174, Brussels, Belgium. Association for Comput-	Guillaume Lample, Lucile Saulnier, L�elio Ren-	764
711	ational Linguistics.	ard Lavaud, Marie-Anne Lachaux, Pierre Stock,	765
712	Yen-Chun Chen and Mohit Bansal. 2018. Fast abstrac-	Teven Le Scao, Thibaut Lavril, Thomas Wang, Timo-	766
713	tive summarization with reinforce-selected sentence	th�ee Lacroix, and William El Sayed. 2023. Mistral	767
714	rewriting. <i>arXiv preprint arXiv:1805.11080</i> .	7b . ArXiv:2310.06825 [cs].	768
715	Sumit Chopra, Michael Auli, and Alexander M Rush.	Shubhra Kanti Karmaker Santu, Chase Geigle, Dun-	769
716	2016. Abstractive sentence summarization with at-	cun Ferguson, William Cope, Mary Kalantzis, Du-	770
717	tentive recurrent neural networks. In <i>Proceedings of</i>	ane Searsmith, and Chengxiang Zhai. 2018. Sofsat:	771
718	<i>the 2016 Conference of the North American Chap-</i>	Towards a setlike operator based framework for se-	772
719	<i>ter of the Association for Computational Linguistics:</i>	mantic analysis of text . <i>ACM SIGKDD Explorations</i>	773
720	<i>Human Language Technologies</i> , pages 93–98.	<i>Newsletter</i> , 20(2):21–30.	774
721	Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Zhifang Sui,	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	775
722	and Furu Wei. 2022. Why can gpt learn in-context?	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gon-	776
723	language models secretly perform gradient descent as	zalez, Hao Zhang, and Ion Stoica. 2023. Efficient	777
724	meta optimizers. <i>arXiv preprint arXiv:2212.10559</i> .	memory management for large language model serv-	778
725	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	ing with pagedattention . In <i>Proceedings of the 29th</i>	779
726	Kristina Toutanova. 2019. BERT: Pre-training of	<i>Symposium on Operating Systems Principles, SOSP</i>	780
727	deep bidirectional transformers for language under-	'23, page 611–626, New York, NY, USA. Association	781
728	standing . In <i>Proceedings of the 2019 Conference of</i>	for Computing Machinery.	782

783	Logan Lebanoff, Kaiqiang Song, and Fei Liu. 2018.	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	837
784	Adapting the neural encoder-decoder framework	Dario Amodei, and Ilya Sutskever. Language models	838
785	from single to multi-document summarization. <i>arXiv</i>	are unsupervised multitask learners.	839
786	<i>preprint arXiv:1808.06218</i> .		
787	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	840
788	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	841
789	Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: De-	Wei Li, and Peter J Liu. 2019. Exploring the limits	842
790	noising sequence-to-sequence pre-training for natural	of transfer learning with a unified text-to-text trans-	843
791	language generation, translation, and comprehension.	former. <i>arXiv preprint arXiv:1910.10683</i> .	844
792	<i>arXiv preprint arXiv:1910.13461</i> .		
793	Chin-Yew Lin. 2004. ROUGE: A package for auto-	Alberto D. Rodriguez, Katherine R. Dearstyne, and Jane	845
794	matic evaluation of summaries. In <i>Text Summariza-</i>	Cleland-Huang. 2023. Prompts matter: Insights and	846
795	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.	strategies for prompt engineering in automated soft-	847
796	Association for Computational Linguistics.	ware traceability. In <i>2023 IEEE 31st International</i>	848
		<i>Requirements Engineering Conference Workshops</i>	849
		(<i>REW</i>), page 455–464.	850
797	Linqing Liu, Yao Lu, Min Yang, Qiang Qu, Jia Zhu, and	Alexander M Rush, Sumit Chopra, and Jason We-	851
798	Hongyan Li. 2017. Generative adversarial network	ston. 2015. A neural attention model for ab-	852
799	for abstractive text summarization. <i>arXiv preprint</i>	stractive sentence summarization. <i>arXiv preprint</i>	853
800	<i>arXiv:1711.09357</i> .	<i>arXiv:1509.00685</i> .	854
801	Congbo Ma, Wei Emma Zhang, Mingyu Guo, Hu Wang,	Victor Sanh, Lysandre Debut, Julien Chaumond, and	855
802	and Quan Z Sheng. 2020. Multi-document sum-	Thomas Wolf. 2019. Distilbert, a distilled version	856
803	marization via deep learning techniques: A survey.	of bert: smaller, faster, cheaper and lighter. <i>ArXiv</i> ,	857
804	<i>arXiv preprint arXiv:2011.04843</i> .	abs/1910.01108.	858
805	Yogesh Kumar Meena, Ashish Jain, and Dinesh	Shubhra Kanti Karmaker Santu and Dongji Feng. 2023.	859
806	Gopalani. 2014. Survey on graph and cluster based	Teler: A general taxonomy of llm prompts for bench-	860
807	approaches in multi-document text summarization.	marking complex tasks. In <i>Findings of the Associa-</i>	861
808	In <i>International Conference on Recent Advances and</i>	<i>tion for Computational Linguistics: EMNLP 2023</i> ,	862
809	<i>Innovations in Engineering (ICRAIE-2014)</i> , pages	page 14197–14203, Singapore. Association for Com-	863
810	1–5. IEEE.	putational Linguistics.	864
811	Ramesh Nallapati, Bowen Zhou, Caglar Gulcehre, Bing	Shubhra Kanti Karmaker Santu, Chase Geigle, Dun-	865
812	Xiang, et al. 2016. Abstractive text summarization	can Ferguson, William Cope, Mary Kalantzis, Du-	866
813	using sequence-to-sequence rnns and beyond. <i>arXiv</i>	ane Searsmith, and Chengxiang Zhai. 2018. Sofsat:	867
814	<i>preprint arXiv:1602.06023</i> .	Towards a setlike operator based framework for se-	868
815	Shashi Narayan, Shay B Cohen, and Mirella Lapata.	matic analysis of text. <i>ACM SIGKDD Explorations</i>	869
816	2018. Ranking sentences for extractive summariza-	<i>Newsletter</i> , 20(2):21–30.	870
817	tion with reinforcement learning. <i>arXiv preprint</i>	Souvika Sarkar, Dongji Feng, and Shubhra Kanti Kar-	871
818	<i>arXiv:1802.08636</i> .	maker Santu. 2022. Exploring universal sentence	872
819	OpenAI. 2023. Gpt-4 technical report.	encoders for zero-shot text classification. In <i>Pro-</i>	873
820	ArXiv:2303.08774 [cs].	<i>ceedings of the 2nd Conference of the Asia-Pacific</i>	874
821	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-	<i>Chapter of the Association for Computational Lin-</i>	875
822	<i>guistics and the 12th International Joint Conference</i>	876	
823	Sandhini Agarwal, Katarina Slama, Alex Ray, John	<i>on Natural Language Processing (Volume 2: Short</i>	877
824	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	<i>Papers)</i> , page 135–147, Online only. Association for	878
825	Maddie Simens, Amanda Askell, Peter Welinder,	Computational Linguistics.	879
826	Paul Christiano, Jan Leike, and Ryan Lowe. 2022.	Souvika Sarkar, Dongji Feng, and Shubhra Kanti Kar-	880
827	Training language models to follow instructions with	maker Santu. 2023. Zero-shot multi-label topic infer-	881
828	human feedback. ArXiv:2203.02155 [cs].	ence with sentence encoders and llms. In <i>Proceed-</i>	882
829	Romain Paulus, Caiming Xiong, and Richard Socher.	<i>ings of the 2023 Conference on Empirical Methods</i>	883
830	2017. A deep reinforced model for abstractive sum-	<i>in Natural Language Processing</i> , page 16218–16233,	884
831	marization. <i>arXiv preprint arXiv:1705.04304</i> .	Singapore. Association for Computational Linguis-	885
832	Alec Radford, Jeffrey Wu, Dario Amodei, Daniela	tics.	886
833	Amodei, Jack Clark, Miles Brundage, and Ilya	Stanford Law School. 2023. Large language models	887
834	Sutskever. 2019. Better language models and	as fiduciaries: A case study toward robustly com-	888
835	their implications. <i>OpenAI Blog</i> https://openai.	municating with artificial intelligence through legal	889
836	com/blog/better-language-models .	standards.	890

891	Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank	narayana Pillai, Jacob Devlin, Michael Laskin, Diego	952
892	Tip. 2024. An empirical evaluation of using large	de Las Casas, Dasha Valter, Connie Tao, Lorenzo	953
893	language models for automated unit test genera-	Blanco, Adrià Puigdomènech Badia, David Reitter,	954
894	tion. <i>IEEE Transactions on Software Engineering,</i>	Mianna Chen, Jenny Brennan, Clara Rivera, Sergey	955
895	50(1):85–105.	Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski,	956
		Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yim-	957
896	Nan Shao, Zefan Cai, Chonghua Liao, Yanan Zheng,	ing Gu, Kate Olszewska, Ravi Addanki, Antoine	958
897	Zhilin Yang, et al. 2023. Compositional task rep-	Miech, Annie Louis, Denis Teplyashin, Geoff Brown,	959
898	resentations for large language models. In <i>The</i>	Elliot Catt, Jan Balaguer, Jackie Xiang, Pidong Wang,	960
899	<i>Eleventh International Conference on Learning Rep-</i>	Zoe Ashwood, Anton Briukhov, Albert Webson, San-	961
900	<i>resentations.</i>	jay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-	962
901	Shaden Smith, Mostofa Patwary, Brandon Norick,	Wei Chang, Axel Stjerngren, Josip Djolonga, Yut-	963
902	Patrick LeGresley, Samyam Rajbhandari, Jared	ing Sun, Ankur Bapna, Matthew Aitchison, Pedram	964
903	Casper, Zhun Liu, Shrimai Prabhumoye, George	Pejman, Henryk Michalewski, Tianhe Yu, Cindy	965
904	Zerveas, Vijay Korthikanti, et al. 2022. Using deep-	Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich,	966
905	speed and megatron to train megatron-turing nlg	Kehang Han, Peter Humphreys, Thibault Sellam,	967
906	530b, a large-scale generative language model. <i>arXiv</i>	James Bradbury, Varun Godbole, Sina Samangooei,	968
907	<i>preprint arXiv:2201.11990.</i>	Bogdan Damoc, Alex Kaskasoli, Sébastien M. R.	969
		Arnold, Vijay Vasudevan, Shubham Agrawal, Jason	970
908	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-	Riesa, Dmitry Lepikhin, Richard Tanburn, Srivat-	971
909	Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan	san Srinivasan, Hyeontaek Lim, Sarah Hodkinson,	972
910	Schalkwyk, Andrew M. Dai, Anja Hauth, Katie	Pranav Shyam, Johan Ferret, Steven Hand, Ankush	973
911	Millican, David Silver, Melvin Johnson, Ioannis	Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Gi-	974
912	Antonoglou, Julian Schrittwieser, Amelia Glaese,	ang, Alexander Neitz, Zaheer Abbas, Sarah York,	975
913	Jilin Chen, Emily Pitler, Timothy Lillicrap, Ange-	Machel Reid, Elizabeth Cole, Aakanksha Chowdh-	976
914	liki Lazaridou, Orhan Firat, James Molloy, Michael	ery, Dipanjan Das, Dominika Rogozińska, Vitaliy	977
915	Isard, Paul R. Barham, Tom Hennigan, Benjamin	Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas	978
916	Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong	Zilka, Flavien Prost, Luheng He, Marianne Mon-	979
917	Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza	teiro, Gaurav Mishra, Chris Welty, Josh Newlan,	980
918	Rutherford, Erica Moreira, Kareem Ayoub, Megha	Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu,	981
919	Goel, Jack Krawczyk, Cosmo Du, Ed Chi, Heng-	Raoul de Liedekerke, Justin Gilmer, Carl Saroufim,	982
920	Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty	Shruti Rijhwani, Shaobo Hou, Disha Shrivastava,	983
921	Chan, Manaal Faruqui, Aliaksei Severyn, Hanzhao	Anirudh Baddepudi, Alex Goldin, Adnan Ozturel,	984
922	Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah,	Albin Cassirer, Yunhan Xu, Daniel Sohn, Deven-	985
923	Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran,	dra Sachan, Reinald Kim Amplayo, Craig Swans-	986
924	Sumit Bagri, Balaji Lakshminarayanan, Jeremiah	son, Dessie Petrova, Shashi Narayan, Arthur Guez,	987
925	Liu, Andras Orban, Fabian Göra, Hao Zhou, Xiny-	Siddhartha Brahma, Jessica Landon, Miteyan Pa-	988
926	ing Song, Aurelien Boffy, Harish Ganapathy, Steven	tel, Ruizhe Zhao, Kevin Vellela, Luyu Wang, Wen-	989
927	Zheng, HyunJeong Choe, Ágoston Weisz, Tao Zhu,	hao Jia, Matthew Rahtz, Mai Giménez, Legg Yeung,	990
928	Yifeng Lu, Siddharth Gopal, Jarrod Kahn, Maciej	James Keeling, Petko Georgiev, Diana Mincu, Boxi	991
929	Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa,	Wu, Salem Haykal, Rachel Saputro, Kiran Vodra-	992
930	Majd Al Mery, Martin Baeuml, Zhifeng Chen, Lau-	halli, James Qin, Zeynep Cankara, Abhanshu Sharma,	993
931	rent El Shafey, Yujing Zhang, Olcan Sercinoglu,	Nick Fernando, Will Hawkins, Behnam Neyshabur,	994
932	George Tucker, Enrique Piqueras, Maxim Krikun,	Solomon Kim, Adrian Hutter, Priyanka Agrawal,	995
933	Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca	Alex Castro-Ros, George van den Driessche, Tao	996
934	Roelofs, Anaïs White, Anders Andreassen, Tamara	Wang, Fan Yang, Shuo-yiin Chang, Paul Komarek,	997
935	von Glehn, Lakshman Yagati, Mehran Kazemi, Lu-	Ross McIlroy, Mario Lučić, Guodong Zhang, Wael	998
936	cas Gonzalez, Misha Khalman, Jakub Sygnowski,	Farhan, Michael Sharman, Paul Natsev, Paul Michel,	999
937	Alexandre Frechette, Charlotte Smith, Laura Culp,	Yamini Bansal, Siyuan Qiao, Kris Cao, Siamak Shak-	1000
938	Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan	eri, Christina Butterfield, Justin Chung, Paul Kishan	1001
939	Schucher, Federico Lebron, Alban Rrustemi, Nati-	Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar	1002
940	alie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao,	Soparkar, Karel Lenc, Timothy Chung, Aedan Pope,	1003
941	Bartek Perz, Dian Yu, Heidi Howard, Adam Blo-	Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo	1004
942	niarz, Jack W. Rae, Han Lu, Laurent Sifre, Mar-	Wang, Joshua Maynez, Mary Phuong, Taylor Tobin,	1005
943	cello Maggioni, Fred Alcober, Dan Garrette, Megan	Andrea Tacchetti, Maja Trebacz, Kevin Robinson,	1006
944	Barnes, Shantanu Thakoor, Jacob Austin, Gabriel	Yash Katariya, Sebastian Riedel, Paige Bailey, Kefan	1007
945	Barth-Maron, William Wong, Rishabh Joshi, Rahma	Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone,	1008
946	Chaabouni, Deeni Fatiha, Arun Ahuja, Gaurav Singh	Neil Houlsby, Xuehan Xiong, Zhen Yang, Elena Gri-	1009
947	Tomar, Evan Senter, Martin Chadwick, Ilya Kor-	bovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music	1010
948	nakov, Nithya Attaluri, Iñaki Iturrate, Ruibo Liu,	Li, Thais Kagohara, Jay Pavagadhi, Sophie Bridgers,	1011
949	Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao Jia,	Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed,	1012
950	Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse	Tianqi Liu, Richard Powell, Vijay Bolina, Mariko	1013
951	Hartman, Xavier Garcia, Thanumalayan Sankara-	inuma, Polina Zablotskaia, James Besley, Da-Woon	1014

1015	Chung, Timothy Dozat, Ramona Comanescu, Xiance Si, Jeremy Greer, Guolong Su, Martin Polacek,	Liu, Jules Walter, Hamid Moghaddam, Arun Kishore,	1078
1016	Raphaël Lopez Kaufman, Simon Tokumine, Hexiang	Jakub Adamek, Tyler Mercado, Jonathan Mallinson,	1079
1017	Hu, Elena Buchatskaya, Yingjie Miao, Mohamed	Siddhinita Wandekar, Stephen Cagle, Eran Ofek,	1080
1018	Elhawaty, Aditya Siddhant, Nenad Tomasev, Jin-	Guillermo Garrido, Clemens Lombriser, Maksim	1081
1019	wei Xing, Christina Greer, Helen Miller, Shereen	Mukha, Botu Sun, Hafeezul Rahman Mohammad,	1082
1020	Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Ange-	Josip Matak, Yadi Qian, Vikas Peswani, Pawel Janus,	1083
1021	los Filos, Milos Besta, Rory Blevins, Ted Klimenko,	Quan Yuan, Leif Schelin, Oana David, Ankur Garg,	1084
1022	Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Os-	Yifan He, Oleksii Duzhyi, Anton Algmyr, Timo-	1085
1023	car Chang, Mantas Pajarskas, Carrie Muir, Vered	thée Lottaz, Qi Li, Vikas Yadav, Luyao Xu, Alex	1086
1024	Cohen, Charline Le Lan, Krishna Haridasan, Amit	Chinien, Rakesh Shivanna, Aleksandr Chuklin, Josie	1087
1025	Marathe, Steven Hansen, Sholto Douglas, Rajku-	Li, Carrie Spadine, Travis Wolfe, Kareem Mohamed,	1088
1026	mar Samuel, Mingqiu Wang, Sophia Austin, Chang	Subhabrata Das, Zihang Dai, Kyle He, Daniel von	1089
1027	Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso Lorenzo,	Dincklage, Shyam Upadhyay, Akanksha Maurya,	1090
1028	Lars Lowe Sjöstrand, Sébastien Cevey, Zach Gle-	Luyan Chi, Sebastian Krause, Khalid Salama, Pam G.	1091
1029	icher, Thi Avrahami, Anudhyan Boral, Hansa Srimi-	Rabinovitch, Pavan Kumar Reddy M, Aarush Sel-	1092
1030	vasan, Vittorio Selo, Rhys May, Konstantinos Aiso-	van, Mikhail Dektiarev, Golnaz Ghiasi, Erdem Gu-	1093
1031	pos, Léonard Hussenot, Livio Baldini Soares, Kate	ven, Himanshu Gupta, Boyi Liu, Deepak Sharma,	1094
1032	Baumli, Michael B. Chang, Adrià Recasens, Ben	Idan Heimlich Shtacher, Shachi Paul, Oscar Aker-	1095
1033	Caine, Alexander Pritzel, Filip Pavetic, Fabio Pardo,	lund, François-Xavier Aubet, Terry Huang, Chen	1096
1034	Anita Gergely, Justin Frye, Vinay Ramasesh, Dan	Zhu, Eric Zhu, Elico Teixeira, Matthew Fritze,	1097
1035	Horgan, Kartikeya Badola, Nora Kassner, Subhra-	Francesco Bertolini, Liana-Eleonora Marinescu, Mar-	1098
1036	jit Roy, Ethan Dyer, Víctor Campos Campos, Alex	tin Bölle, Dominik Paulus, Khyatti Gupta, Tejasi	1099
1037	Tomala, Yunhao Tang, Dalia El Badawy, Elspeth	Latkar, Max Chang, Jason Sanders, Roopa Wil-	1100
1038	White, Basil Mustafa, Oran Lang, Abhishek Jin-	son, Xuwei Wu, Yi-Xuan Tan, Lam Nguyen Thiet,	1101
1039	dal, Sharad Vikram, Zhitao Gong, Sergi Caelles,	Tulsee Doshi, Sid Lall, Swaroop Mishra, Wanming	1102
1040	Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng,	Chen, Thang Luong, Seth Benjamin, Jasmine Lee,	1103
1041	Wojciech Stokowiec, Ce Zheng, Phoebe Thacker,	Ewa Andrejczuk, Dominik Rabiej, Vipul Ranjan,	1104
1042	Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh,	Krzysztof Styrac, Pengcheng Yin, Jon Simon, Mal-	1105
1043	James Svensson, Max Bileschi, Piyush Patil, Ankesh	colm Rose Harriott, Mudit Bansal, Alexei Robsky,	1106
1044	Anand, Roman Ring, Katerina Tsihlias, Arpi Vezer,	Geoff Bacon, David Greene, Daniil Mirylenka, Chen	1107
1045	Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom	Zhou, Obaid Sarvana, Abhimanyu Goyal, Samuel	1108
1046	Kwiatkowski, Samira Daruki, Keran Rong, Allan	Andermatt, Patrick Siegler, Ben Horn, Assaf Is-	1109
1047	Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg,	rael, Francesco Pongetti, Chih-Wei “Louis” Chen,	1110
1048	Mina Khan, Lisa Anne Hendricks, Marie Pellat,	Marco Selvatici, Pedro Silva, Kathie Wang, Jack-	1111
1049	Vladimir Feinberg, James Cobon-Kerr, Tara Sainath,	son Tolins, Kelvin Guu, Roey Yogev, Xiaochen Cai,	1112
1050	Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives,	Alessandro Agostini, Maulik Shah, Hung Nguyen,	1113
1051	Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd,	Noah Ó Donnaile, Sébastien Pereira, Linda Friso,	1114
1052	Le Hou, Qingze Wang, Thibault Sottiaux, Michela	Adam Stambler, Adam Kurzrok, Chenkai Kuang,	1115
1053	Paganini, Jean-Baptiste Lespiau, Alexandre Mou-	Yan Romanikhin, Mark Geller, Z. J. Yan, Kane Jang,	1116
1054	farek, Samer Hassan, Kaushik Shivakumar, Joost van	Cheng-Chun Lee, Wojciech Fica, Eric Malmi, Qi-	1117
1055	Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh	jun Tan, Dan Banica, Daniel Balle, Ryan Pham,	1118
1056	Goyal, Matthew Tung, Andrew Brock, Hannah Sheah-	Yanping Huang, Diana Avram, Hongzhi Shi, Jasjot	1119
1057	an, Vedant Misra, Cheng Li, Nemanja Rakićević,	Singh, Chris Hidey, Niharika Ahuja, Pranab Saxe-	1120
1058	Mostafa Dehghani, Fangyu Liu, Sid Mittal, Jun-	na, Dan Dooley, Srividya Pranavi Potharaju, Eileen	1121
1059	hyuk Oh, Seb Noury, Eren Sezener, Fantine Huot,	O’Neill, Anand Gokulchandran, Ryan Foley, Kai	1122
1060	Matthew Lamm, Nicola De Cao, Charlie Chen, Sid-	Zhao, Mike Dusenberry, Yuan Liu, Pulkit Mehta,	1123
1061	harth Mudgal, Romina Stella, Kevin Brooks, Gau-	Ragha Kotikalapudi, Chalence Safranek-Shrader, An-	1124
1062	tam Vasudevan, Chenxi Liu, Mainak Chain, Nivedita	drew Goodman, Joshua Kessinger, Eran Globen, Pra-	1125
1063	Melinkeri, Aaron Cohen, Venus Wang, Kristie Sey-	teek Kolhar, Chris Gorgolewski, Ali Ibrahim, Yang	1126
1064	more, Sergey Zubkov, Rahul Goel, Summer Yue,	Song, Ali Eichenbaum, Thomas Brovelli, Sahitya	1127
1065	Sai Krishnakumaran, Brian Albert, Nate Hurley,	Potluri, Preethi Lahoti, Cip Baetu, Ali Ghorbani,	1128
1066	Motoki Sano, Anhad Mohanany, Jonah Joughin,	Charles Chen, Andy Crawford, Shalini Pal, Mukund	1129
1067	Egor Filonov, Tomasz Kępa, Yomna Eldawy, Jiaw-	Sridhar, Petru Gurita, Asier Mujika, Igor Petrovski,	1130
1068	ern Lim, Rahul Rishi, Shirin Badiezedegan, Taylor	Pierre-Louis Cedoz, Chenmei Li, Shiyuan Chen,	1131
1069	Bos, Jerry Chang, Sanil Jain, Sri Gayatri Sundara	Nicolò Dal Santo, Siddharth Goyal, Jitesh Pun-	1132
1070	Padmanabhan, Subha Puttagunta, Kalpesh Krishna,	jabi, Karthik Kappaganthu, Chester Kwak, Pallavi	1133
1071	Leslie Baker, Norbert Kalb, Vamsi Bedapudi, Adam	LV, Sarmishta Velury, Himadri Choudhury, Jamie	1134
1072	Kurzrok, Shuntong Lei, Anthony Yu, Oren Litvin,	Hall, Premal Shah, Ricardo Figueira, Matt Thomas,	1135
1073	Xiang Zhou, Zhichun Wu, Sam Sobell, Andrea Si-	Minjie Lu, Ting Zhou, Chintu Kumar, Thomas Ju-	1136
1074	ciliano, Alan Papir, Robby Neale, Jonas Bragagnolo,	rdis, Sharat Chikkerur, Yenai Ma, Adams Yu, Soo	1137
1075	Tej Toor, Tina Chen, Valentin Anklin, Feiran Wang,	Kwak, Victor Åhdel, Sujevan Rajayogam, Travis	1138
1076	Richie Feng, Milad Gholami, Kevin Ling, Lijuan	Choma, Fei Liu, Aditya Barua, Colin Ji, Ji Ho	1139
1077		Park, Vincent Hellendoorn, Alex Bailey, Taylan Bi-	1140

1141	lal, Huanjie Zhou, Mehrdad Khatir, Charles Sutton, Wojciech Rzadkowski, Fiona Macintosh, Konstantin Shagin, Paul Medina, Chen Liang, Jinjing Zhou, Pararth Shah, Yingying Bi, Attila Dankovics, Shipra Banga, Sabine Lehmann, Marissa Bredesen, Zifan Lin, John Eric Hoffmann, Jonathan Lai, Raynald Chung, Kai Yang, Nihal Balani, Arthur Bražinskas, Andrei Sozanschi, Matthew Hayes, Héctor Fernández Alcalde, Peter Makarov, Will Chen, Antonio Stella, Liselotte Snijders, Michael Mandl, Ante Kärman, Paweł Nowak, Xinyi Wu, Alex Dyck, Krishnan Vaidyanathan, Raghavender R, Jessica Mallet, Mitch Rudominer, Eric Johnston, Sushil Mittal, Akhil Udathu, Janara Christensen, Vishal Verma, Zach Irving, Andreas Santucci, Gamaleldin Elsayed, Elnaz Davoodi, Marin Georgiev, Ian Tenney, Nan Hua, Geoffrey Cideron, Edouard Leurent, Mahmoud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy Zheng, Dylan Scandinaro, Heinrich Jiang, Jasper Snoek, Mukund Sundararajan, Xuezhi Wang, Zack Ontiveros, Itay Karo, Jeremy Cole, Vinu Rajashekhar, Lara Tumeh, Eyal Ben-David, Rishub Jain, Jonathan Uesato, Romina Datta, Oskar Bunyan, Shimu Wu, John Zhang, Piotr Stanczyk, Ye Zhang, David Steiner, Subhajit Naskar, Michael Azzam, Matthew Johnson, Adam Paszke, Chung-Cheng Chiu, Jaume Sanchez Elias, Afroz Mohiuddin, Faizan Muhammad, Jin Miao, Andrew Lee, Nino Vieillard, Jane Park, Jiageng Zhang, Jeff Stanway, Drew Garmon, Abhijit Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Luowei Zhou, Jonathan Evens, William Isaac, Geoffrey Irving, Edward Loper, Michael Fink, Isha Arkatkar, Nanxin Chen, Izhak Shafran, Ivan Petrychenko, Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai Zhu, Peter Grabowski, Yu Mao, Alberto Magni, Kaisheng Yao, Javier Snaider, Norman Casagrande, Evan Palmer, Paul Suganthan, Alfonso Castaño, Irene Giannoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy Wiesner, Diana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian LIN, Marcus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Ginger Perng, Elena Allica Abellan, Mingyang Zhang, Ishita Dasgupta, Nate Kushman, Ivo Penchev, Alena Repina, Xihui Wu, Tom van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse, Fan Yang, Jeff Piper, Nathan Ie, Rama Pasmarthi, Nathan Lintz, Anitha Vijayakumar, Daniel Andor, Pedro Valenzuela, Minnie Lui, Cosmin Padurar, Daiyi Peng, Katherine Lee, Shuyuan Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhal, Dayou Du, Dan McKinnon, Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Ken Franko, Anna Bulanova, Rémi Leblond, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu,	
	Felix Fischer, Jun Xu, Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev, Hannah Forbes, Dylan Banarse, Zora Tung, Mark Omernick, Colton Bishop, Rachel Sterneck, Rohan Jain, Jiawei Xia, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Daniel J. Mankowitz, Alex Polozov, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu, Meghana Thotakuri, Tom Natan, Matthieu Geist, Ser tan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Danila Sinopalnikov, Sabela Ramos, John Mellor, Abhishek Sharma, Kathy Wu, David Miller, Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang, Rui Zhu, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Daniel Toyama, Evan Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina, Rishika Sinha, Alice Talbert, Diane Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana, Jing Li, Saaber Fatehi, John Wieting, Omar Ajmeri, Benigno Urias, Yeongil Ko, Laura Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Yeqing Li, Nir Levine, Ariel Stolovich, Rebecca Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Charlie Deck, Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem, Sho Arora, Christy Koh, Soheil Hassas Yeganeh, Siim Pöder, Mukarram Tariq, Yanhua Sun, Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, T. J. Lu, Abe Ittycheriah, Prakash Shroff, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David Gaddy, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza Ghaffarkhah, Morgane Rivière, Alanna Walton, Clément Crepy, Alicia Parrish, Zongwei Zhou, Clement Farabet, Carey Radebaugh, Praveen Srinivasan, Claudia van der Salm, Andreas Fildjeland, Salvatore Scellato, Eri Latorre-Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Cesare, Tom Hudson, Piermaria Mendolichio, Lexi Walker, Alex Morris, Matthew Mauger, Alexey Guseynov, Alison Reid, Seth Odoom, Lucia Loher, Victor Cotruta, Madhavi Yenugula, Dominik Grewe, Anastasia Petrushkina, Tom Duerig, Antonio Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson, Lynette Webb, Sahil Dua, Dong Li, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj Khare, Shreyas Rammohan Belle, Lei Wang, Chetan Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao Zhao, Stephan	1204 1205 1206 1207 1208 1209 1210 1211 1212 1213 1214 1215 1216 1217 1218 1219 1220 1221 1222 1223 1224 1225 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237 1238 1239 1240 1241 1242 1243 1244 1245 1246 1247 1248 1249 1250 1251 1252 1253 1254 1255 1256 1257 1258 1259 1260 1261 1262 1263 1264 1265 1266

1267	Lee, Pandu Nayak, Doug Fritz, Manish Reddy	new standard for open-source, commercially usable	1329
1268	Vuyyuru, John Aslanides, Nidhi Vyas, Martin Wicke,	llms . Accessed: 2024-01-30.	1330
1269	Xiao Ma, Evgenii Eltyshv, Nina Martin, Hardie		
1270	Cate, James Manyika, Keyvan Amiri, Yelin Kim,		
1271	Xi Xiong, Kai Kang, Florian Luisier, Nilesh Tripu-	Arun James Thirunavukarasu, Darren Shu Jeng Ting,	1331
1272	raneni, David Madras, Mandy Guo, Austin Waters,	Kabilan Elangovan, Laura Gutierrez, Ting Fang	1332
1273	Oliver Wang, Joshua Ainslie, Jason Baldridge, Han	Tan, and Daniel Shu Wei Ting. 2023. Large	1333
1274	Zhang, Garima Pruthi, Jakob Bauer, Feng Yang,	language models in medicine . <i>Nature Medicine</i> ,	1334
1275	Riham Mansour, Jason Gelman, Yang Xu, George	29(88):1930–1940.	1335
1276	Polovets, Ji Liu, Honglong Cai, Warren Chen, Xiang-	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	1336
1277	Hai Sheng, Emily Xue, Sherjil Ozair, Christof Anger-	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	1337
1278	mueller, Xiaowei Li, Anoop Sinha, Weiren Wang, Ju-	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	1338
1279	lia Wiesinger, Emmanouil Koukoumidis, Yuan Tian,	Azhar, Aurelien Rodriguez, Armand Joulin, Edouard	1339
1280	Anand Iyer, Madhu Gurusurthy, Mark Goldenson,	Grave, and Guillaume Lample. 2023a. Llama:	1340
1281	Parashar Shah, M. K. Blake, Hongkun Yu, Anthony	Open and efficient foundation language models .	1341
1282	Urbanowicz, Jennimaria Palomaki, Chrisantha Fer-	ArXiv:2302.13971 [cs].	1342
1283	nando, Ken Durden, Harsh Mehta, Nikola Mom-	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	1343
1284	chev, Elahe Rahimtoroghi, Maria Georgaki, Amit	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	1344
1285	Raul, Sebastian Ruder, Morgan Redshaw, Jinhyuk	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	1345
1286	Lee, Denny Zhou, Komal Jalan, Dinghua Li, Blake	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	1346
1287	Hechtman, Parker Schuh, Milad Nasr, Kieran Milan,	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	1347
1288	Vladimir Mikulik, Juliana Franco, Tim Green, Nam	Jude Fernandes, Jeremy Fu, Wenxin Fu, Brian Fuller,	1348
1289	Nguyen, Joe Kelley, Aroma Mahendru, Andrea Hu,	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	1349
1290	Joshua Howland, Ben Vargas, Jeffrey Hui, Kshitij	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	1350
1291	Bansal, Vikram Rao, Rakesh Ghiya, Emma Wang,	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	1351
1292	Ke Ye, Jean Michel Sarr, Melanie Moranski Preston,	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	1352
1293	Madeleine Elish, Steve Li, Aakash Kaku, Jigar Gupta,	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	1353
1294	Ice Pasupat, Da-Cheng Juan, Milan Someswar, Tejvi	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	1354
1295	M., Xinyun Chen, Aida Amini, Alex Fabrikant, Eric	tinnet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	1355
1296	Chu, Xuanyi Dong, Amruta Muthal, Senaka Buth-	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	1356
1297	pitiya, Sarthak Jauhari, Nan Hua, Urvashi Khan-	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	1357
1298	delwal, Ayal Hitron, Jie Ren, Larissa Rinaldi, Sha-	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	1358
1299	har Drath, Avigail Dabush, Nan-Jiang Jiang, Har-	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	1359
1300	shal Godhia, Uli Sachs, Anthony Chen, Yicheng	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	1360
1301	Fan, Hagai Taitelbaum, Hila Noga, Zhuyun Dai,	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	1361
1302	James Wang, Chen Liang, Jenny Hamer, Chun-Sung	Melanie Kambadur, Sharan Narang, Aurelien Ro-	1362
1303	Ferng, Chenel Elkind, Aviel Atias, Paulina Lee, Vít	driguez, Robert Stojnic, Sergey Edunov, and Thomas	1363
1304	Listfk, Mathias Carlen, Jan van de Kerkhof, Marcin	Scialom. 2023b. Llama 2: Open foundation and	1364
1305	Pikus, Krunoslav Zaher, Paul Müller, Sasha Zykova,	fine-tuned chat models . ArXiv:2307.09288 [cs].	1365
1306	Richard Stefanec, Vitaly Gatsko, Christoph Hirn-	Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt	1366
1307	schall, Ashwin Sethi, Xingyu Federico Xu, Chetan	Haberland, Tyler Reddy, David Cournapeau, Ev-	1367
1308	Ahuja, Beth Tsai, Anca Stefanoiu, Bo Feng, Ke-	geni Burovski, Pearu Peterson, Warren Weckesser,	1368
1309	shav Dhandhania, Manish Katyal, Akshay Gupta,	Jonathan Bright, Stéfan J. van der Walt, Matthew	1369
1310	Atharva Parulekar, Divya Pitta, Jing Zhao, Vivaan	Brett, Joshua Wilson, K. Jarrod Millman, Nikolay	1370
1311	Bhatia, Yashodha Bhavnani, Omar Alhadlaq, Xiaolin	Mayorov, Andrew R. J. Nelson, Eric Jones, Robert	1371
1312	Li, Peter Danenberg, Dennis Tu, Alex Pine, Vera	Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng,	1372
1313	Filippova, Abhipso Ghosh, Ben Limonchik, Bhar-	Eric W. Moore, Jake VanderPlas, Denis Laxalde,	1373
1314	gava Urala, Chaitanya Krishna Lanka, Derik Clive,	Josef Perktold, Robert Cimrman, Ian Henriksen, E. A.	1374
1315	Yi Sun, Edward Li, Hao Wu, Kevin Hongtongsak,	Quintero, Charles R. Harris, Anne M. Archibald, An-	1375
1316	Ianna Li, Kalind Thakkar, Kuanysh Omarov, Kushal	tônio H. Ribeiro, Fabian Pedregosa, Paul van Mul-	1376
1317	Majmundar, Michael Alverson, Michael Kucharski,	bregt, and SciPy 1.0 Contributors. 2020. SciPy 1.0:	1377
1318	Mohak Patel, Mudit Jain, Maksim Zabelin, Paolo	Fundamental Algorithms for Scientific Computing in	1378
1319	Pelagatti, Rohan Kohli, Saurabh Kumar, Joseph Kim,	Python . <i>Nature Methods</i> , 17:261–272.	1379
1320	Swetha Sankar, Vineet Shah, Lakshmi Ramachan-	Shomir Wilson, Florian Schaub, Aswarth Abhilash	1380
1321	dru, Xiangkai Zeng, Ben Bariach, Laura Weidinger,	Dara, Frederick Liu, Sushain Cherivirala, Pedro	1381
1322	Tu Vu, Amar Subramanya, Sissie Hsiao, Demis Hass-	Giovanni Leon, Mads Schaarup Andersen, Sebas-	1382
1323	abis, Koray Kavukcuoglu, Adam Sadovsky, Quoc	tian Zimmeck, Kanthashree Mysore Sathyendra,	1383
1324	Le, Trevor Strohman, Yonghui Wu, Slav Petrov,	N. Cameron Russell, Thomas B. Norton, Eduard	1384
1325	Jeffrey Dean, and Oriol Vinyals. 2024. Gemini:	Hovy, Joel Reidenberg, and Norman Sadeh. 2016.	1385
1326	A family of highly capable multimodal models .	The creation and analysis of a website privacy policy	1386
1327	(arXiv:2312.11805). ArXiv:2312.11805 [cs].	corpus . In <i>Proceedings of the 54th Annual Meet-</i>	1387
1328	MosaicML NLP Team. 2023. Introducing mpt-7b: A	ing of the Association for Computational Linguistics	1388

1389	(<i>Volume 1: Long Papers</i>), page 1330–1340, Berlin, Germany. Association for Computational Linguistics.	1445
1390		1446
1391	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Huggingface’s transformers: State-of-the-art natural language processing . ArXiv:1910.03771 [cs].	1447
1392		1448
1393		1449
1394		1450
1395		1451
1396		1452
1397		1453
1398		1454
1399		
1400	Yuxiang Wu and Baotian Hu. 2018. Learning to extract coherent summary via deep reinforcement learning. <i>arXiv preprint arXiv:1804.07036</i> .	
1401		
1402		
1403	Dongling Xiao, Han Zhang, Yukun Li, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. 2020. Ernie-gen: An enhanced multi-flow pre-training and fine-tuning framework for natural language generation. <i>arXiv preprint arXiv:2001.11314</i> .	
1404		
1405		
1406		
1407		
1408	Yu Yan, Weizhen Qi, Yeyun Gong, Dayiheng Liu, Nan Duan, Jiusheng Chen, Ruofei Zhang, and Ming Zhou. 2020. Prophetnet: Predicting future n-gram for sequence-to-sequence pre-training. <i>arXiv preprint arXiv:2001.04063</i> .	
1409		
1410		
1411		
1412		
1413	Michihiro Yasunaga, Rui Zhang, Kshitijh Meelu, Ayush Pareek, Krishnan Srinivasan, and Dragomir Radev. 2017. Graph-based neural multi-document summarization. <i>arXiv preprint arXiv:1706.06681</i> .	
1414		
1415		
1416		
1417	Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J Liu. 2019. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. <i>arXiv preprint arXiv:1912.08777</i> .	
1418		
1419		
1420		
1421	Jinming Zhao, Ming Liu, Longxiang Gao, Yuan Jin, Lan Du, He Zhao, He Zhang, and Gholamreza Haffari. 2020. Summpip: Unsupervised multi-document summarization with sentence graph compression. In <i>Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 1949–1952.	
1422		
1423		
1424		
1425		
1426		
1427		
1428	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhonghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena . ArXiv:2306.05685 [cs].	
1429		
1430		
1431		
1432		
1433		
1434	Ming Zhong, Pengfei Liu, Yiran Chen, Danqing Wang, Xipeng Qiu, and Xuanjing Huang. 2020. Extractive summarization as text matching. <i>arXiv preprint arXiv:2004.08795</i> .	
1435		
1436		
1437		
1438	Ming Zhong, Pengfei Liu, Danqing Wang, Xipeng Qiu, and Xuanjing Huang. 2019. Searching for effective neural extractive summarization: What works and what’s next. <i>arXiv preprint arXiv:1907.03491</i> .	
1439		
1440		
1441		
1442	Qingyu Zhou, Nan Yang, Furu Wei, and Ming Zhou. 2017. Selective encoding for abstractive sentence summarization. <i>arXiv preprint arXiv:1704.07073</i> .	
1443		
1444		
	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2022. Large language models are human-level prompt engineers. <i>arXiv preprint arXiv:2211.01910</i> .	1445
		1446
		1447
		1448
	Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. 2021. A robustly optimized bert pre-training approach with post-training . In <i>Proceedings of the 20th Chinese National Conference on Computational Linguistics</i> , page 1218–1227, Huhhot, China. Chinese Information Processing Society of China.	1449
		1450
		1451
		1452
		1453
		1454
	A Appendix	1455
	A.1 Additional Figures	1456
	Figure 4 shows a comprehensive breakdown of the best scores obtained by each model for each dataset. Figure 5 shows Pearson’s correlation scores between all metrics on both datasets. The Pearson scores were computed using the SciPy library (Virtanen et al., 2020)	1457
		1458
		1459
		1460
		1461
		1462
	A.2 More on Prompt Design	1463
	The prompt designs for each group mostly follow the format covered in section 3.2 but the entire taxonomy is best laid out by Figure 6. The exact prompts are laid out in the following passage.	1464
		1465
		1466
	System Role Templates Our system role templates are made up of 2 AllSides-specific items, 2 3P specific-items and 6 for general purpose. These are written as follows	1468
		1469
		1470
		1471
	• AllSides	1472
	– you will be given two news articles to read. then you will be given an instruction. follow these instructions as closely as possible	1473
		1474
		1475
	– you will read two news articles and answer any questions about them	1476
		1477
	• 3P	1478
	– you are to read two privacy policies and briefly provide information according to the user’s needs	1479
		1480
	– you are to read two privacy policies and provide concise answers to the user	1481
		1482
	• Both	1483
	– you are to read several documents and briefly provide information according to the user’s needs	1484
		1485
	– you are to read several documents and provide concise answers to the user	1486
		1487
	– you will read two documents and give brief answers to user questions	1488
		1489
	– you are a machine who is given 3 inputs: document 1, document 2, and the instructions. your output will adhere to these 3 inputs.	1490
		1491
		1492
	– you will be given 2 documents and a set of instructions. follow the instructions as closely as possible.	1493
		1494
	– you will be given 2 documents and a set of instructions. your response to these instructions will rely on the material covered in the 2 documents.	1495
		1496
		1497
	In-Context Learning Template: We use the following for our in-context learning template:	1498
		1499

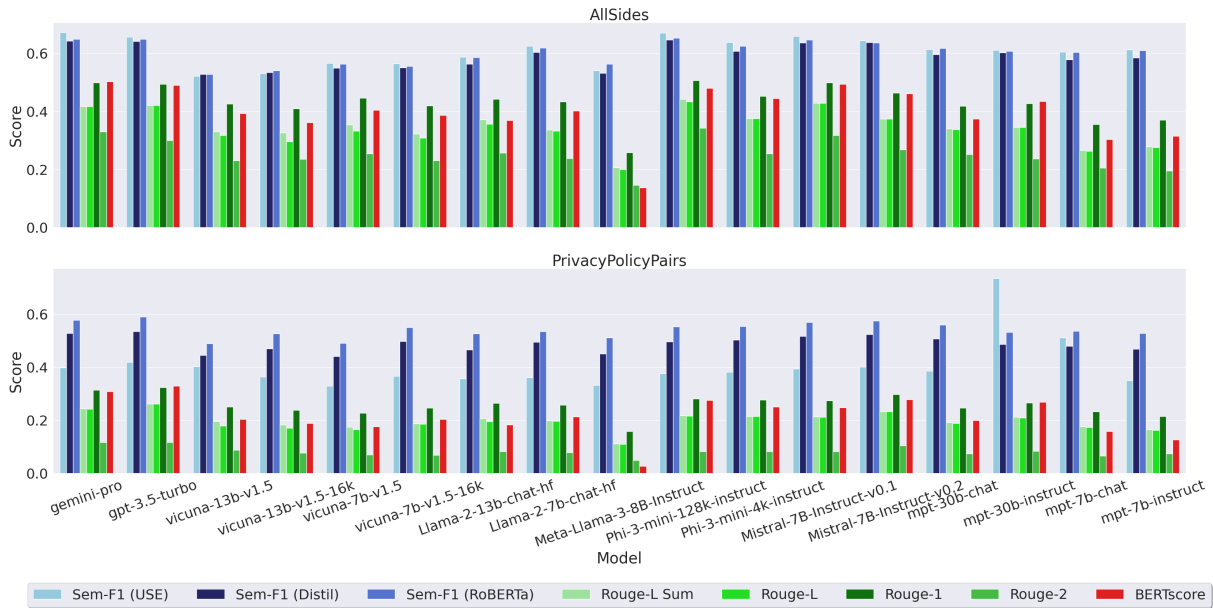


Figure 4: Best scores over each TeLeR prompt level for all 16 evaluated LLMs and for each dataset. Red shows BERTscore, green shows ROUGE, and blue shows Sem-F1.

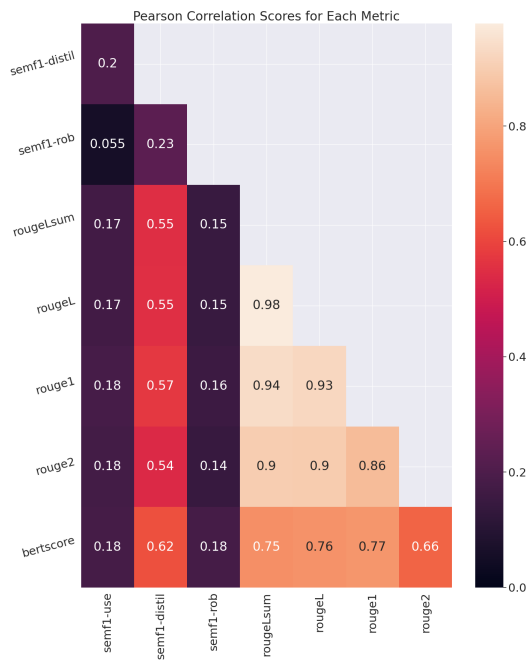


Figure 5: Correlation scores between all evaluation metrics.

- Document 1: **{{Example Document 1}}**
 Document 2: **{{Example Document 2}}**
 Summary: **{{Example Reference}}**
- Document 1: **{{Document 1}}**
 Document 2: **{{Document 2}}**
 Summary:

TeLeR Level 0 Template: With no possibility for variation, our TeLeR L0 template is written as follows:

- {Document 1} {Document 2}** 1510
- TeLeR Level 1 Template:** For our TeLeR L1 templates we have 3 AllSides-only items, 3 3P-only items, and 5 general-purpose items. 1511-1513
- AllSides** 1514
 - Document 1: **{{Document 1}}** 1515
 Document 2: **{{Document 2}}** 1516-1517
 In one sentence, please tell me the overlapping information between article 1 and article 2 1518-1519
 - Document 1: **{{Document 1}}** 1520
 Document 2: **{{Document 2}}** 1521-1522
 summarize the overlapping information between the articles 1523-1524
 - Document 1: **{{Document 1}}** 1525
 Document 2: **{{Document 2}}** 1526-1527
 output the overlapping information of the events covered in these articles 1528-1529
- 3P** 1530
 - Policy 1: **{{Document 1}}** 1531
 Policy 2: **{{Document 2}}** 1532-1533
 In one sentence, please tell me the overlapping information between policy 1 and policy 2 1534-1535
 - Policy 1: **{{Document 1}}** 1536
 Policy 2: **{{Document 2}}** 1537-1538
 summarize the information that the two policies share 1539
 - Policy 1: **{{Document 1}}** 1540
 Policy 2: **{{Document 2}}** 1541-1542
 what is the shared information between the two policies 1543-1544
- Both** 1545

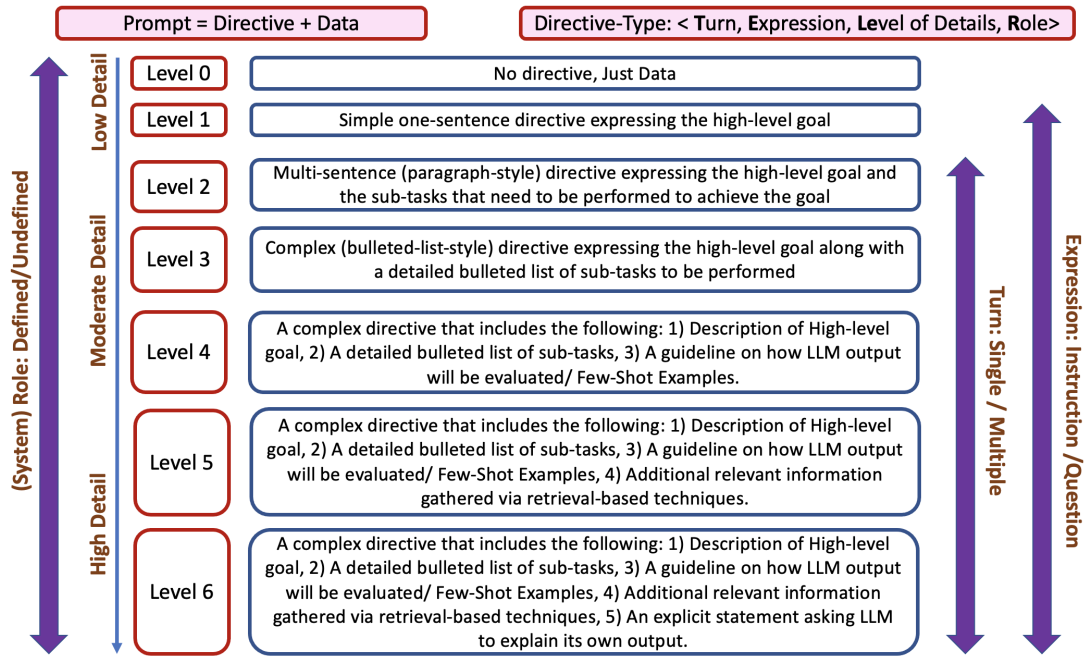


Figure 6: TELeR Taxonomy proposed by Santu and Feng (2023): (<Turn, Expression, Level of Details, Role>)

1546	- Document 1: {{Document 1}}	- Document 1: {{Document 1}}	1582
1547	Document 2: {{Document 2}}	Document 2: {{Document 2}}	1583
1548			1584
1549	In one sentence, please tell me the overlapping information between Document 1 and Document 2	who or what are the common subjects of the two documents? what events are common between the documents? do the documents mention any locations that are the same between the two? give your response in a single sentence.	1585
1550			1586
1551	- Document 1: {{Document 1}}		1587
1552	Document 2: {{Document 2}}		1588
1553			1589
1554	summarize the overlapping information between the documents.	- Document 1: {{Document 1}}	1590
1555		Document 2: {{Document 2}}	1591
1556	- Document 1: {{Document 1}}		1592
1557	Document 2: {{Document 2}}	summarize the overlap	1593
1558			1594
1559	output the overlapping information between the documents.	• 3P	1595
1560		- Policy 1: {{Document 1}}	1596
1561	- Document 1: {{Document 1}}	Policy 2: {{Document 2}}	1597
1562	Document 2: {{Document 2}}		1598
1563		These policies are categorized under "Category". Describe the common aspects of these two policies in terms of this category. make sure to include the shared entities, actions and scope of the documents. Do not make any mention of information that is not shared between them. Keep your response short	1599
1564	output the common information between the documents.		1600
1565			1601
1566	- Document 1: {{Document 1}}		1602
1567	Document 2: {{Document 2}}		1603
1568			1604
1569	output only the overlapping information	- Policy 1: {{Document 1}}	1605
		Policy 2: {{Document 2}}	1606
1570	TELeR Level 2 Templates: For our TELeR L2 templates we have 3 AllSides-only items, 3 3P-only items, and 3 general-purpose items.	These policies are categorized under "Category". Describe the common aspects of these two policies in terms of this category. make sure to include the shared entities, actions and scope of the documents. Do not make any mention of information that is not shared between them. give your response in a single sentence.	1607
1571			1608
1572			1609
1573	• AllSides		1610
1574	- Document 1: {{Document 1}}		1611
1575	Document 2: {{Document 2}}		1612
1576			1613
1577	these articles share similarities. output the information that is shared between them. keep your output short. to be as accurate as possible, cover the "who, what, when, where, and why of the shared information.	- Policy 1: {{Document 1}}	1614
1578		Policy 2: {{Document 2}}	1615
1579			1616
1580		These privacy policy excerpts are tagged with the category: "Category". summarize the overlapping information between the documents. to be as accurate	1617
1581			1618
			1619

1620	as possible, cover the who, what, when, where, and	These policies are labelled under the "Category" cat-	1684
1621	why of the common information.	egory. With this in mind, use a single sentence that	1685
1622	• Both	answers the following:	1686
1623	– Document 1: {{Document 1}}	- Describe the common aspects of these two policies	1687
1624	Document 2: {{Document 2}}	in terms of this category.	1688
1625		- make sure to include the shared entities, actions and	1689
1626	summarize the overlapping information between the	scope of the documents.	1690
1627	two documents. explain the who, what, when, where,	- Do not make any mention of information that is not	1691
1628	and why to give full context.	shared between them.	1692
1629	– Document 1: {{Document 1}}	- Do not respond in a list format and instead respond	1693
1630	Document 2: {{Document 2}}	normally.	1694
1631		• Policy 1: {{Document 1}}	1695
1632	summarize the overlapping information between the	Policy 2: {{Document 2}}	1696
1633	two documents. explain the who, what, when, where,		1697
1634	and why to give full context. the output should be	These policies are labelled under the "Category" cat-	1698
1635	two sentences at most.	egory. With this in mind, use a single sentence that	1699
1636	– Document 1: {{Document 1}}	answers the following:	1700
1637	Document 2: {{Document 2}}	- summarize the information that is shared between	1701
1638		the policies	1702
1639	output the shared information between the documents.	- cover the who, what, when, where, and why of the	1703
1640	do not include any information outside of the shared	common information	1704
1641	information. keep your response short.	- respond in as few sentences as possible	1705
		• Both	1706
1642	TELeR Level 3 Templates: For our TELeR L3	– Document 1: {{Document 1}}	1707
1643	templates we have 3 AllSides-only items, 3 3P-only	Document 2: {{Document 2}}	1708
1644	items, and 2 general-purpose items.		1709
1645	• AllSides	please answer the following:	1710
1646	– Document 1: {{Document 1}}	- who or what are the common subjects of the two	1711
1647	Document 2: {{Document 2}}	documents	1712
1648		- what events are common between the documents	1713
1649	please answer the following:	- do the documents mention any locations that are the	1714
1650	- who or what are the common subjects of the two	same between the two	1715
1651	documents	- keep your response brief. 2 sentences max.	1716
1652	- what events are common between the documents	– Document 1: {{Document 1}}	1717
1653	- do the documents mention any locations that are the	Document 2: {{Document 2}}	1718
1654	same between the two		1719
1655	- keep your response brief. 2 sentences max.	Consider the following questions and respond in a	1720
1656	– Document 1: {{Document 1}}	single sentence:	1721
1657	Document 2: {{Document 2}}	- who or what are the common subjects of the two	1722
1658		documents	1723
1659	Consider the following questions and respond in a	- what events are common between the documents	1724
1660	single sentence:	- do the documents mention any locations that are the	1725
1661	- who or what are the common subjects of the two	same between the two	1726
1662	documents		
1663	- what events are common between the documents	TELeR Level 4 Templates For our TELeR L4	1727
1664	- do the documents mention any locations that are the	templates we have 3 AllSides-only items, 3 3P-	1728
1665	same between the two	only items, and 2 general-purpose items.	1729
1666	• 3P		
1667	– Policy 1: {{Document 1}}	• AllSides	1730
1668	Policy 2: {{Document 2}}	– Document 1: {{Document 1}}	1731
1669		Document 2: {{Document 2}}	1732
1670	These policies are categorized under "Category".		1733
1671	With this in mind, please answer the following:	your goal is to describe all the common information	1734
1672	- Describe the common aspects of these two policies	between the given documents. to accomplish this you	1735
1673	in terms of this category.	will need to answer the following:	1736
1674	- make sure to include the shared entities, actions and	- who or what are the common subjects of the two	1737
1675	scope of the documents.	documents	1738
1676	- Do not make any mention of information that is not	- what events are common between the documents	1739
1677	shared between them.	- do the documents mention any locations that are the	1740
1678	- Do not respond in a list format and instead respond	same between the two	1741
1679	normally.	- keep your response brief. 2 sentences max.	1742
1680	- Keep your response to 3 sentences at most		1743
1681	– Policy 1: {{Document 1}}	For Example:	1744
1682	Policy 2: {{Document 2}}	Doc1: i have a dog. it's pretty fast.	1745
1683		Doc2: i have a dog. he is a slow runner	1746
		Reference Summary: i have a dog.	1747

1748	- Document 1: {{Document 1}}	Reference Summary: Both sentences talk about the	1816
1749	Document 2: {{Document 2}}	speed of a dog	1817
1750		- Policy 1: {{Document 1}}	1818
1751	your goal is to describe all the common information	Policy 2: {{Document 2}}	1819
1752	between the given documents. to accomplish this you		1820
1753	will need to answer the following:	your goal is to describe all the common information	1821
1754	- who or what are the common subjects of the two	between the given documents in one sentence. your	1822
1755	documents	single-sentence response will need to include the	1823
1756	- what events are common between the documents	following:	1824
1757	- do the documents mention any locations that are the	- common aspects related to the given category	1825
1758	same between the two	- common entities	1826
1759		- common applications	1827
1760	your response will be evaluated according to how		1828
1761	similar it is to a "reference summary".	your response will be evaluated according to how	1829
1762	Example:	similar it is to a "reference summary".	1830
1763	Question: what is common between the sentence "the		1831
1764	dog is slow" and "the dog is fast"	Example Documents:	1832
1765	Reference Summary: Both sentences talk about the	Doc1: the dog is slow	1833
1766	speed of a dog	Doc2: the dog is fast	1834
1767			1835
1768	- Document 1: {{Document 1}}	Example Response:	1836
1769	Document 2: {{Document 2}}	Both sentences talk about the speed of a dog	1837
1770	your goal is to describe all the common information	- Policy 1: {{Document 1}}	1838
1771	between the given documents in one sentence. your	Policy 2: {{Document 2}}	1839
1772	single-sentence response will need to capture the		1840
1773	following:	your goal is to describe all the common information	1841
1774	- the common events	between the given documents in one sentence. your	1842
1775	- common people	single-sentence response will need to include the	1843
1776	- common locations	following:	1844
1777	- the overlapping narrative of the documents	- common aspects related to the given category	1845
1778		- common entities	1846
1779	your response will be evaluated according to how	- common applications	1847
1780	similar it is to a "reference summary".		1848
1781	Example:	your response will be evaluated according to how	1849
1782	Doc1: the dog is slow	similar it is to a "reference summary".	1850
1783	Doc2: the dog is fast		1851
1784	Reference Summary: Both sentences talk about the	Example Documents:	1852
1785	speed of a dog	Doc1: the dog is slow	1853
1786		Doc2: the dog is fast	1854
1787	• 3P		1855
1788	- Policy 1: {{Document 1}}	Example Response:	1856
1789	Policy 2: {{Document 2}}	Both sentences talk about the speed of a dog	1857
1790	your goal is to describe all the common information		1858
1791	between the given privacy policies. to accomplish	• Both	1859
1792	this you will need to answer according to the	- Document 1: {{Document 1}}	1860
1793	following:	Document 2: {{Document 2}}	1861
1794	- Describe the common aspects of these two policies	Write a summary of the given documents that follows	1862
1795	in terms of this category.	these instructions:	1863
1796	- make sure to include the shared entities, actions and	- who or what are the common subjects of the two	1864
1797	scope of the documents.	documents	1865
1798	- Do not make any mention of information that is not	- what events are common between the documents	1866
1799	shared between them.	- do the documents mention any locations that are the	1867
1800	- Do not respond in a list format and instead respond	same between the two	1868
1801	normally.	- keep your response brief. 2 sentences max.	1869
1802	- Keep your response to 3 sentences at most		1870
1803		your response will be evaluated according to how	1871
1804	your response will be evaluated according to how	similar it is to a "reference summary".	1872
1805	similar it is to a "reference summary".	For Example:	1873
1806	For example, an output of "cat" could be compared to	Doc1: i have a dog. it's pretty fast.	1874
1807	"light" to get a score of 0 but that same output could	Doc2: i have a dog. he is a slow runner	1875
1808	be compared to "cat" to receive a score of 100. These	Reference Summary: i have a dog.	1876
1809	reference summaries are usually quite short so it is		1877
1810	important to keep your response to 3 sentences or less.	- Document 1: {{Document 1}}	1878
1811		Document 2: {{Document 2}}	1879
1812	your response will be evaluated according to how	Summarize the overlapping information between	1880
1813	similar it is to a "reference summary". Example:	these documents. your summary should follow these	1881
1814	Doc1: the dog is slow	instructions:	1882
1815	Doc2: the dog is fast	- exclude any information that is similar but differing	1883

1884 or contradictory
1885 - write the summary as if you were summarizing a
1886 single document.
1887 - your summary should be short. keep it within 2
1888 sentences.
1889
1890 your response will be evaluated according to how
1891 similar it is to a "reference summary".
1892 For Example:
1893 Doc1: i have a dog. it's pretty fast.
1894 Doc2: i have a dog. he is a slow runner
1895 Reference Summary: i have a dog.