ACADEMICEVAL: LIVE LONG-CONTEXT LLM BENCHMARK

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) have achieved remarkable performance in longcontext understanding. However, current long-context LLM benchmarks are limited by rigid context length and labor-intensive annotation, and the label leakage issue in LLM training also poses a pressing challenge. Therefore, we propose ACADEMICEVAL, a live benchmark for evaluating LLMs over long-context generation tasks. ACADEMICEVAL adopts papers on arXiv to introduce several academic writing tasks with long-context inputs, *i.e.*, TITLE, ABSTRACT, INTRO-DUCTION, and RELATED WORK, which cover a wide range of abstraction levels and require no manual labeling. Moreover, ACADEMICEVAL integrates highquality and expert-curated few-shot demonstrations from a collected co-author graph to enable flexible context length. Especially, ACADEMICEVAL features an efficient live evaluation, ensuring no label leakage. We conduct holistic experiments on ACADEMICEVAL, and the results illustrate that LLMs perform poorly on tasks with hierarchical abstraction levels and tend to struggle with long fewshot demonstrations, illustrating the challenge of our benchmark. We also provide insightful analysis for enhancing LLMs' long-context modeling capabilities.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 029

Large Language Models (LLMs) have recently achieved tremendous success in natural language
processing (NLP) tasks Achiam et al. (2023); AI@Meta (2024). However, when facing long context
inputs, LLMs show a sharp decline in performance, which poses a pressing challenge to LLMs
in understanding and capturing key information in long texts Li et al. (2024); Liu et al. (2024).
Therefore, several long-context LLM benchmarks are spawned to evaluate LLMs in various settings,
including question answering, summarizing, and reasoning Shaham et al. (2023); An et al. (2023);
Dong et al. (2023); Bai et al. (2023b); Li et al. (2023); Zhang et al. (2024). Despite their success,
these benchmarks still suffer from concerns of rigid context length, saturated performance, and being
leaked in LLM training.

038 We envision that the *next-generation long-context LLM benchmarks* should ideally possess three key features. (1) Flexible and potentially unlimited context length: existing benchmarks fix the context 040 for each long-context problem; ideally, the format and length of the context could be flexibly set 041 based on the LLM's capability, especially given the release of long-context LLMs Reid et al. (2024) 042 and their capabilities in ingesting multi-modal information, e.g., graphs Dong et al. (2024). (2) 043 High-quality labels derived from *real-world data*, *minimizing* human labeling efforts: existing long-044 context benchmarks often require human labeling Bai et al. (2023b); An et al. (2023); Li et al. (2023); Dong et al. (2023); Zhang et al. (2024), which is costly and limits the size of the benchmarks to about 2000 samples Xu et al. (2023) (3) Live updates to mitigate information leakage during LLM 046 pretraining and fine-tuning: benchmark data contamination in LLM has gradually become a severe 047 issue Sainz et al. (2023); Ye et al. (2024); Zhu et al. (2024b;a); Xu et al. (2024); we argue that holding 048 out future data as the val/test set is one of the most effective approaches for open benchmarks. 049

Based on these principles, we propose ACADEMICEVAL, a live benchmark to evaluate LLMs over
 long-context generation tasks. ACADEMICEVAL adopts arXiv as its data source and features a suite
 of academic writing tasks on each paper without labor-intensive annotation: TITLE, ABSTRACT,
 INTRODUCTION, and RELATED WORK, each of which has long-context input and hierarchical abstraction levels. In particular, we construct a co-author graph via arXiv API to conveniently obtain

054

056

058

Table 1: Comparison with Existing Long-context LLM Benchmarks. Each column indicates the
average input length, whether the annotation is human-assisted, whether there are tasks with hierar-
chical abstraction levels, whether it contains few-shot demonstrations, and whether the benchmark
is lively updated, respectively.

Benchmark	Avg Len	Automatic Annotation	Hierarchical Abstraction	Few-shot Demons	Live Update
ZeroSCROLLS (Shaham et al., 2023)	$\sim 10 \text{K}$	 Image: A set of the set of the	×	×	×
L-Eval (An et al., 2023)	$\sim 8 \mathrm{K}$	×	×	×	×
BAMBOO (Dong et al., 2023)	~16K	×	×	×	×
LongBench (Bai et al., 2023b)	$\sim 8 K$	×	×	1	×
LooGLE (Li et al., 2023)	$\sim 20 K$	×	×	×	×
∞ Bench (Zhang et al., 2024)	$\sim 200 \mathrm{K}$	×	×	×	×
AcademicEval (ours)	Flexible	✓	1	✓	1

067 068 069

co-author papers as high-quality and expert-curated few-shot demonstrations, which also possess ACADEMICEVAL flexible context length. Furthermore, ACADEMICEVAL introduces efficient live 071 evaluation based on the co-author graph, which utilizes the latest papers on arXiv to update the 072 benchmark data periodically and ensures no label leakage. Moreover, ACADEMICEVAL provides 073 in-context few-shot demonstrations for each data sample, which is neglected by most existing long-074 context LLM benchmarks Liu et al. (2024); Li et al. (2024). In the experiment, we evaluate three 075 types of LLMs on ACADEMICEVAL: standard LLMs, long-context LLMs, and retrieval-augmented 076 language models (RALM). Experimental results show that current LLMs cannot deal with long-077 context context tasks well at diverse abstraction levels, and RALM is the worst-performing one among the three types of baselines. Additionally, as the input length increases, noticeable perfor-079 mance degradation can be seen on almost all tasks, with the largest drop reaching 32% and 7% w.r.t. RougeL Lin (2004) and BERTScore Zhang et al. (2019), respectively. Although we find that fewshot demonstrations from co-author papers can slightly strengthen the performance over some tasks, 081 it is still limited by the long context modeling capabilities of LLMs. In general, the experimental findings indicate that ACADEMICEVAL is a challenging long-context LLM benchmark. 083

We illustrate the comparison with existing long-context LLM benchmarks in Table 1. Our contribu-tions are summarized as follows:

- We propose a live benchmark, ACADEMICEVAL, to evaluate LLMs over long-context generation tasks. ACADEMICEVAL features four academic writing tasks with hierarchical abstraction levels and requires no manual annotation.
 - We construct a co-author graph via the arXiv API and draw on the co-author papers as informative few-shot demonstrations, making the context length of ACADEMICEVAL flexible and scalable.
 - ACADEMICEVAL conducts periodic data updates on the co-author graph to enable efficient live evaluation, which ensures no label leakage and fair evaluation.
 - We conduct comprehensive experiments on ACADEMICEVAL, demonstrating its challenges and providing insights for improving LLMs in long-context modeling.
- 096 097 098

099

087

090

092

094

095

2 RELATED WORK

Long-context Modeling and LLM Benchmarks LLMs are known to be powerful in language 101 modeling tasks Achiam et al. (2023); AI@Meta (2024). However, when it comes to long-context 102 input, LLMs show a sharp decline in performance, posing a pressing challenge when benchmark-103 ing their long-context modeling capabilities Liu et al. (2024); Li et al. (2024). Currently, there are 104 two mainstream technologies for long-context modeling tasks: retrieval-augmented language mod-105 els (RALM)Ram et al. (2023); Yu et al. (2023); Trivedi et al. (2022); Jiang et al. (2023); Asai et al. (2023) and long-context LLMs Bai et al. (2023a); Jiang et al. (2024); Teknium et al.. RALM equips 106 LLMs with a retrieverRobertson et al. (2009); Ramos et al. (2003); Karpukhin et al. (2020); Izac-107 ard et al. (2021) to perform information retrieval on short text chunks, which are then fed to LLMs together with the input query to generate the final output. As a retrieval system, RALM is usually
evaluated over retrieval-based benchmarks, including STARK Wu et al. (2024), RGB Chen et al.
(2024), ARES Saad-Falcon et al. (2023), etc. In comparison, long-context LLMs expand their context window length to accommodate longer inputs and are benchmarked over various tasks, which
include long-context QA, summarization, conversations, reasoning, etc Shaham et al. (2023); An
et al. (2023); Dong et al. (2023); Bai et al. (2023b); Li et al. (2023); Zhang et al. (2024).

114 Label Leakage in LLM Benchmarks Label leakage has always been a severe issue that bench-115 marks must attempt to avoid during data collection. However, recent researches Xu et al. (2024); 116 Zhu et al. (2024b;a); Ye et al. (2024) point out that most LLM benchmarks are composed of statically 117 collected data, which may be inevitably included in the large amount of training data of LLMs, caus-118 ing label leakage. Therefore, some works attempt to measure or detect the extent of label leakage in LLM benchmarks. Benbench Xu et al. (2024) leverages perplexity and N-gram accuracy to quantify 119 potential label leakage, while PAC Ye et al. (2024) detects contaminated data by comparing the po-120 larized distance of samples before and after augmentation. Even though these approaches propose 121 to measure or detect label leakage, there is little work on mitigating and solving this issue Zhu et al. 122 (2024b). Dynabench Kiela et al. (2021) and Dynaboard Ma et al. (2021) feature dynamic human-123 in-the-loop dataset creation while avoiding leakage, which is very labor-intensive. DyVal Zhu et al. 124 (2024b) leverages pre-set constraints and directed acyclic graphs (DAG) to dynamically generate test 125 cases with diverse complexities, reducing the risk of label leakage. FreshBench Zhu et al. (2024a) 126 and StackMIA Ye et al. (2024) collect the latest data from public websites periodically and simply 127 rely on the chronological split to build a dynamic benchmark. 128

Long-context Summarization Benchmarks Solving ACADEMICEVAL requires LLM's long-129 context summarization capability Liu et al. (2024). Existing works include (1) query-based summa-130 rization tasks, focusing on the capability of models to position and capture local key information in 131 long texts given a specific query Litvak & Vanetik (2017); Wang et al. (2022); (2) single-document 132 or multi-document summarization tasks concentrate on evaluating the ability of models to under-133 stand long texts holistically Cohan et al. (2018); Meng et al. (2021); Huang et al. (2021); Kryściński 134 et al. (2021); Cachola et al. (2020). These long-context summarization benchmarks suffer from the 135 above-mentioned limitations, including requiring human-assisted labeling and concerns about data leakage; moreover, these summarization tasks focus on one-level summarization, failing to consider 136 the summarizations at different abstraction levels. 137

138 139

140 141

3 ACADEMICEVAL BENCHMARK

In this section, we propose ACADEMICEVAL (Figure 1) for live evaluation in long-context generation tasks with hierarchical abstraction levels. We first describe data collection and preprocessing in Section 3.1. Then, in Section 3.2, four academic writing tasks with diverse abstraction levels are introduced, and we also integrate few-shot demonstrations to make the context length flexible and scalable. Finally, Section 3.3 elucidates the live evaluation with periodic data updates.

146 147 148

3.1 DATA CURATION

149 **Co-author Graph Construction via arXiv** As a public paper preprint platform, arXiv¹ has always 150 been favored by researchers. It archives a huge amount of papers and updates the latest ones daily, 151 which serves as an excellent data source and also lays the foundation for the live update of our 152 benchmark. Thanks to the arXiv API², paper files can be obtained in batch without much manual 153 effort. We first collect and construct a co-author graph (edges are established between two author 154 nodes that are co-authors in at least one paper) using the arXiv API through breadth-first search 155 (BFS), where the features of each author node include the first-author papers published by the author. 156 By making the co-author graph the carrier of paper data, we can form an interconnected whole of scattered papers, which provides valuable structural information to be exploited for our benchmark 157 (e.g., as few-shot demonstrations). Furthermore, we can enable efficient live updates on the co-158 author graph, which will be introduced in Section 3.3. 159

¹⁶⁰ 161

¹https://arxiv.org/

²https://info.arxiv.org/help/api/index.html



Figure 1: AcademicEval Benchmark. We construct a co-author graph via arXiv and conduct a 178 chronological split on all paper samples (training, validation, and test samples are represented by 179 red, orange, and green, respectively). Each paper sample is preprocessed into separate sections and can be integrated with few-shot demonstrations from co-author papers.

- 181 182
- 183

187

191

197

199

205

206

207

208

210

211

213

Academic Data Gathering and Preprocessing After some basic operations on the co-author graph, such as taking the maximum connected component, we preprocess all papers in the node features. 185 For each paper, we collect related metadata via the arXiv API, including author information, publication timestamp, etc., and download the pdf file simultaneously. Intuitively, the paper content can be considered as a kind of original, expert-curated, and high-quality labeled data without manual 188 annotation. Therefore, we develop a complete pipeline for preprocessing each paper, splitting and 189 extracting the text of several sections in it. To give an example, we can utilize the pipeline to extract 190 the introduction section from the main body of a paper. Then, the extracted introduction section, the remaining parts of the extracted main body (i.e., w/o the introduction section), and the abstract and 192 title constitute the basic data of each paper sample. For the related work section, we extract each cited paper's abstract and title to form an additional citation corpus. 193

194 We will describe in detail in Section 3.2 how to use these data to design long-context academic 195 writing tasks. 196

BENCHMARKING LLMS OVER LONG-CONTEXT GENERATION TASKS WITH 3.2 HIERARCHICAL ABSTRACTION

Task Description Employing machine learning approaches to automate academic writing has al-200 ways been a research hotspot with significant practical application value Chen et al. (2022; 2021). 201 Therefore, inspired by the leave-one-out validation, we introduce four academic writing tasks with 202 ultra-long context to evaluate the generation capability of LLMs under different abstraction levels, 203 as shown below: 204

- TITLE WRITING. This task takes a paper's main body and abstract, along with a specific task prompt as inputs, and then asks LLMs to output a predicted title.
- ABSTRACT WRITING. Similar to the above, this task takes a paper's main body and title, along with a specific task prompt as inputs, and then asks LLMs to output a predicted abstract.
- INTRODUCTION WRITING. This task takes a paper's main body (with the introduction section removed), title, and abstract, along with a specific task prompt as inputs, and then 212 asks LLMs to output a predicted introduction.
- RELATED WORK WRITING. This task takes a paper's main body (with the related work 214 section removed), title, abstract, and citation corpus (introduced in Section 3.1), along with 215 a specific task prompt as inputs, and then asks LLMs to output a predicted related work.

Based on the above task descriptions, we can generate four basic benchmark settings with different abstraction levels, namely TITLE-10K, ABS-9K, INTRO-8K and RELATED-34K, with suffixes indicating their input context length³.

219 Integration of Few-shot Demonstrations Given the rigid context length of current long-context 220 LLM benchmarks and the general effectiveness of in-context learning in LLMs Dong et al. (2022); 221 Wei et al. (2022a;b); Kojima et al. (2022), we propose to integrate long few-shot demonstrations 222 to enable flexible and scalable context length, and we have two selection options for each sample 223 in the above four basic benchmark settings: (1) Randomly select papers under the same category. 224 According to the paper categories provided by the arXiv API, we can randomly select several non-225 duplicate papers under the same category. (2) Randomly Select co-author papers. The motivation is 226 straightforward: the similarity of research directions between co-author papers is more fine-grained. Thanks to the co-author graph, it is convenient to obtain co-author papers of each original paper 227 sample. These selected papers serve as few-shot demonstrations and are utilized as input-output 228 pairs to enrich the input context of the original samples, providing potentially insightful and relevant 229 content while enabling flexible and scalable context length. 230

Consequently, we have completed the construction of benchmark settings, and the data statistics inthe initial collection round are shown in Table 2.

233 Data Statistics As shown in Table 2, ACADEMICEVAL has four academic writing tasks with hier-234 archical abstraction levels, and each task features four settings with diverse input context lengths, 235 some of which are obtained by integrating few-shot demonstrations. For instance, each sample in 236 TITLE-10K consists of a single paper sample. TITLE-30K and TITLE-31K-G are obtained by 237 integrating with two few-shot demonstrations from random papers and co-author papers, respec-238 tively, while TITLE-50K-M is obtained by using both of the above integration options. Actually, 239 we can scale context length by increasing the number of few-shot demonstrations to provide more informative references, enhancing task performance. 240

Furthermore, we present the text compression rate (defined as the number of input tokens divided by the number of output tokens) for each benchmark setting in Table 2 to illustrate the diverse abstraction levels in ACADEMICEVAL. Across the four tasks, a higher compression rate means a higher level of text abstraction in this task. Among several settings within each task, a higher compression rate makes it tougher to exploit information holistically but more likely to produce better outputs (since more references are integrated). These different tasks and settings increase the diversity of the ACADEMICEVAL benchmark.

As for data splitting, we perform a chronological split in ACADEMICEVAL, which means that the test set always contains the latest papers collected in each collection round, ensuring no label leakage. Note that Table 2 shows only the data collected in the initial round, which will be updated periodically as described in the next section.

- 252
- 253 254

3.3 LIVE EVALUATION WITH PERIODIC DATA UPDATES ON THE CO-AUTHOR GRAPH

The daily updates of arXiv provide the basis for the live evaluation of ACADEMICEVAL: we can periodically update the benchmark with the latest papers on arXiv. By setting a reasonable update cycle (*e.g.*, monthly or quarterly), we can ensure that the data in the benchmark is not contaminated so that it can be used to evaluate LLMs fairly in a live manner. Therefore, we proposed an efficient incremental update procedure on the co-author graph:

(1) Node Update For each author on the co-author graph, check whether the author has a newly published first-author paper through the arXiv API. If so, add it to the corresponding node feature on the co-author graph.

(2) Node and Edge Update During the traversal of Node Update, each author's new co-authors are added to a candidate list, and the number of new papers (including first-author and non-first-author papers) when searching for the author is used as the priority of the co-authors (co-authors of active authors tend to be active as well, and we can efficiently collect the latest papers from active authors).

 $^{^{3}}$ We use BERT Devlin et al. (2018) tokenizer by default to count the number of input tokens (output tokens are not included).

Table 2: Data Statistics of AcademicEval (Initial Round). It includes 4 writing tasks and provides four settings of different context length for each task. For each setting, we list their Comp. Rate, Samples of Each, Chronological Split, and Timespan of Test Data.

Setting	Comp. Rate (In-Len. / Out-Len.)	#Samples of Each.	Chronological Split (Train-Val-Test)	Timespan of Test Data		
TITLE WRITING						
TITLE-10K TITLE-30K TITLE-31K-G TITLE-50K-M	587 1773 1807 2968	5098	72%-19%-9%	2024.06- 2024.07		
ABSTRACT WRITING						
ABS-9K ABS-28K ABS-29K-G ABS-48K-M	36 108 112 185	5098	72%-19%-9%	2024.06- 2024.07		
INTRODUCTION WRITING						
INTRO-8K INTRO-28K INTRO-28K-G INTRO-48K-M	6 21 22 37	4665	71%-20%-9%	2024.06- 2024.07		
	RELATED V	VORK WRIT	TING			
RELATED-34K RELATED-53K RELATED-53K-G RELATED-72K-M	34 53 53 72	2240	72%-20%-8%	2024.06- 2024.07		

Note: We use the BERT tokenizer by default to count the number of tokens.

Then, we use the prioritized candidate list to conduct BFS to update nodes and edges until a specific number of incremental update papers is met.

(3) Graph Pruning As the benchmark is updated, we will remove some outdated papers and inactive authors (defined as those who have not published new first-author or non-first-author papers for a long time) from the co-author graph.

In this way, the latest papers can be obtained sufficiently and efficiently while ensuring connectivity and a smaller graph size.

Live Leaderboard We also provide a leaderboard for live evaluation of the current most advanced LLMs, which will be released later.

4 EXPERIMENTS

4.1 BASELINES

³¹⁸ We adopt the following three types of baselines to conduct a holistic evaluation of ACADEMICEVAL.

Standard LLMs We choose Gemma Instruct (7B) Team et al. (2024) and LLaMA-3 Chat (70B) AI@Meta (2024) as standard LLM baselines, each with a context length of 8K.

Long-context LLMs We choose Qwen 1.5 Chat (72B) Bai et al. (2023a), Mixtral-8x7B Instruct (46.7B) Jiang et al. (2024), and Nous Hermes 2 - Mixtral 8x7B-DPO (46.7B) Teknium et al. as long-context LLM baselines, each with a context length of 32K.

324 **Retrieval-augmented language models (RALM)** First, we consider two sparse retrievers: (1) 325 BM25 Robertson et al. (2009): This is a widely used retrieval model that ranks documents based 326 on the frequency of query terms in each document. (2) TF-IDF Ramos et al. (2003): It scores doc-327 uments by multiplying the term frequency of each query term by the inverse document frequency. Second, we also consider three dense retrievers: (3) DPR Karpukhin et al. (2020): It uses a bi-328 encoder to retrieve relevant documents based on dense embeddings. (4) Contriever Izacard et al. 329 (2021): It leverages unsupervised contrastive learning to learn high-quality dense representations. 330 (5) Dragon Lin et al. (2023): It enhances retriever training by employing data augmentation, includ-331 ing query and label augmentation. 332

We use the inputs of ACADEMICEVAL as the external corpus of RALM. For text split, we use the RecursiveCharacterTextSplitter from LangChain⁴ and set chunk size and chunk overlap to 512 and 64, respectively. For each retrieval, we recall up to 12 text chunks (limited by the context length of standard LLMs) based on text similarity (semantic similarity based on inner product for dense retrievers or similarity based on word frequency for sparse retrievers).

338 339

4.2 EVALUATION METRICS

For evaluation metrics, we adopt (1) BERTScore⁵ Zhang et al. (2019): This metric leverages
BERT-based embedding to measure semantic similarity between predicted and reference texts. (2)
RougeL Lin (2004): This metric evaluates the longest common subsequence between the generated and reference texts, providing a measure of similarity in terms of sequential matching. For both metrics, higher scores indicate a better match between the predicted and the reference text.

3463474.3 MAIN RESULT ANALYSIS

We conduct comprehensive experiments on the four academic writing tasks, and the results w.r.t. BERTScore and RougeL are presented in Table 3 and 4, respectively. Note that we do not conduct experiments on -M settings because its context length is too long for most of our selected baselines.

351 Diverse Task Difficulties and Abstractions The four tasks we proposed are designed to challenge 352 LLMs over long-context generation tasks with different abstraction levels. From Table 3 and 4, 353 we can clearly observe that it provides different difficulties for LLMs to perform well from TITLE 354 WRITING to RELATED WORK WRITING tasks, and the results of all baselines on these four tasks 355 have a relatively obvious trend. For example, the TITLE WRITING task tends to have a higher score than the ABSTRACT WRITING task, which may indicate that the TITLE WRITING task is easier 356 than the ABSTRACT WRITING task. Since a title only has a few words, LLMs only need to generate 357 a roughly related theme to achieve a high semantic similarity, while an abstract requires a more 358 detailed description generated to achieve it. 359

360 Baseline Performance Comparison Among different baselines, RALM with LLaMA frequently 361 delivers the highest scores across various tasks and context lengths, with only a context length of 8K. Standard LLMs also achieve competitive performance, which is slightly inferior to long-context 362 LLMs. This exposes the shortcomings of long-context LLMs' generation capabilities, which is 363 well revealed by ACADEMICEVAL. Among long-context LLMs, Hermes performs best overall, 364 but is still slightly inferior to RALM with LLaMA. This shows that although the current long-365 context LLMs have a longer context window size, they still have great deficiencies in processing 366 long text information. In contrast, RALM-based methods generally outperform other baselines. This 367 is primarily due to the retrieval mechanisms of RALM, which retrieves and processes information 368 in a few relevant shorter chunks, enabling it to focus on key information. 369

Impact of Context Length The impact of context length on performance is evident across all task settings and both metrics, with baselines generally performing worse as the context length increases. For example, the TITLE WRITING task shows a noticeable drop in scores as the context length extends from 10K to 31K tokens. This trend is also apparent in ABSTRACT WRITING and INTRO-DUCTION WRITING, where longer contexts correlate with decreased model performance. showing that our benchmark challenges LLMs in effectively processing ultra-long inputs.

⁴https://www.langchain.com/

376

⁵We use deberta-xlarge-mnli He et al. (2021) instead of the default roberta-large Liu et al. (2019) as the backbone model to have the best correlation with human evaluation.

M. J.L.	Standard LLMs		Lon	Long-context LLMs			RALM	
Models	Gemma	LLaMA	Qwen	Mixtral	Hermes	Gemma [†]	$LLaMA^{\dagger}$	
#Params.	7B	70B	72B	8x7B	8x7B	7B	70B	
Context Length	8K	8K	32K	32K	32K	8K	8K	
Setting: TITLE W	RITING							
TITLE-10K	66.1	74.1	73.9	73.4	74.2	65.8	73.9	
TITLE-30K	-	-	73.0	72.9	73.4	65.7	73.9	
TITLE-31K-G	-	-	72.8	72.8	73.3	65.7	73.8	
Setting: ABSTRAC	T WRITIN	G						
Авѕ-9К	59.9	62.4	62.5	61.4	62.2	60.3	61.5	
Abs-28K	-	-	61.3	61.2	62.6	60.1	61.4	
Abs-29K-G	-	-	61.3	61.4	62.5	60.2	61.3	
Setting: INTRODUCTION WRITING								
Intro-8K	54.8	55.8	55.4	54.6	55.2	55.0	55.2	
INTRO-28K	-	-	54.8	54.0	54.8	55.0	55.2	
INTRO-28K-G	-	-	54.9	54.1	54.7	55.0	55.3	
Setting: RELATED WORK WRITING								
RELATED-34K	52.0	56.2	58.5	55.3	57.8	52.4	54.7	
Related-53K	-	-	-	-	-	52.4	54.7	
RELATED-53K-G	-	-	-	-	-	52.4	54.8	

Table 3: Main Results on AcademicEval w.r.t. BERTScore.

Bold indicates the highest score in each row.

† denotes augmentation with a retriever (Default: Contriever).

"-" means that the context length is too long to be fed into LLMs.



Figure 2: Analysis of RALM on ABS-9K. The left figure shows results with Gemma Instruct (7B), while the right one shows results with LLaMA-3 Chat (70B).

Impact of Few-shot Demonstrations From Table 3 and 4, we can observe that the integration of few-shot demonstrations generally degrades the performance of baselines, except for a few tasks where the results are slightly improved. This shows that current LLMs cannot exploit long few-shot demonstrations to benefit the target tasks well, emphasizing the importance of evaluating long in-context learning in LLM benchmarks. In addition, we can also find that few-shot demonstrations from co-author papers generally have a more positive impact on task performance than randomly selected ones.

Table 4: Main Results on AcademicEval w.r.t. RougeL.							
Madala	Standard LLMs Long-context LLMs		RALM				
Models	Gemma	LLaMA	Qwen	Mixtral	Hermes	Gemma [†]	LLaMA
#Params.	7B	70B	72B	8x7B	8x7B	7B	70B
Context Length	8K	8K	32K	32K	32K	8K	8K
Setting: TITLE W	RITING						
TITLE-10K	44.5	47.1	44.2	45.2	46.2	42.7	47.3
TITLE-30K	-	-	42.9	44.6	45.9	42.6	47.3
TITLE-31K-G	-	-	44.2	44.4	45.3	42.5	47.0
Setting: ABSTRAC	t Writin	G					
Abs-9K	22.4	25.0	24.3	24.1	26.1	23.4	24.2
Abs-28K	-	-	23.3	24.7	26.6	23.1	24.1
Abs-29K-G	-	-	23.3	24.9	26.6	23.2	24.0
Setting: INTRODU	CTION WI	RITING					
INTRO-8K	14.9	18.1	16.2	17.2	17.8	15.4	17.9
INTRO-28K	-	-	16.3	17.5	17.5	15.3	17.8
INTRO-28K-G	-	-	16.3	17.5	17.5	15.4	17.8
Setting: RELATED	WORK W	RITING					
Related-34K	13.5	14.9	16.0	13.4	15.1	14.1	15.3
Related-53K	-	-	-	-	-	14.0	15.3
Related-53K-G	-	-	-	-	-	14.0	15.2

Bold indicates the highest score in each row.

† denotes augmentation with a retriever (Default: Contriever).

"-" means that the context length is too long to be fed into LLMs.

4.4 ADDITIONAL ANALYSIS ON RALM

We conduct extensive experiments on RALM on the ABS-9K setting using standard LLMs Gemma Instruct (7B) and LLaMA-3 Chat (70B), and the results are presented in Figure 2. We can find that the performance of dense retrievers consistently outperforms sparse retrievers, among which contriever achieves the best results. This is because the summary generation task emphasizes semantic similarity, which can be well measured by the similarity of dense embeddings. However, the sparse retrievers perform text chunk recall based on sparse embeddings, and the results are significantly worse than those of the dense retrievers.

471 472

457

458

463

432

5 CONCLUSION

473 474 475

In this paper, we propose ACADEMICEVAL, a live long-context LLM benchmark for evaluating 476 long-context generation tasks with hierarchical abstraction levels. ACADEMICEVAL adopts arXiv 477 as the data source and introduces several long-context academic writing tasks without manual anno-478 tation since the papers on arXiv can be regarded as original, high-quality, and expert-curated labels. 479 Moreover, we integrate few-shot demonstrations from a collected co-author graph to make the con-480 text length of our benchmark flexible and scalable. An efficient live evaluation is also designed to 481 make ACADEMICEVAL immune to the label leakage issue and move toward a more fair evaluation. 482 In the experiments, we conduct a comprehensive analysis on ACADEMICEVAL using several LLM 483 baselines, and the results show that ACADEMICEVAL is a challenging long-context LLM benchmark. Insightful findings are also elucidated for potentially strengthening the long-context modeling 484 capabilities of LLMs and inspiring future long-context LLM benchmarks to evaluate LLMs more 485 flexibly and holistically.

486 REFERENCES

502

531

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 491 AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/
 492 llama3/blob/main/MODEL_CARD.md.
- Chenxin An, Shansan Gong, Ming Zhong, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu.
 L-eval: Instituting standardized evaluation for long context language models. *arXiv preprint arXiv:2307.11088*, 2023.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023a.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, et al. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*, 2023b.
- Isabel Cachola, Kyle Lo, Arman Cohan, and Daniel S Weld. Tldr: Extreme summarization of scientific documents. *arXiv preprint arXiv:2004.15011*, 2020.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. Benchmarking large language models in
 retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
 volume 38, pp. 17754–17762, 2024.
- Xiuying Chen, Hind Alamro, Mingzhe Li, Shen Gao, Xiangliang Zhang, Dongyan Zhao, and Rui
 Yan. Capturing relations between scientific papers: An abstractive model for related work section
 generation. Association for Computational Linguistics, 2021.
- Xiuying Chen, Hind Alamro, Mingzhe Li, Shen Gao, Rui Yan, Xin Gao, and Xiangliang Zhang.
 Target-aware abstractive related work generation with contrastive learning. In *Proceedings of the 45th international ACM SIGIR conference on research and development in information retrieval*,
 pp. 373–383, 2022.
- Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. A discourse-aware attention model for abstractive summarization of long documents. *arXiv preprint arXiv:1804.05685*, 2018.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Jialin Dong, Bahare Fatemi, Bryan Perozzi, Lin F Yang, and Anton Tsitsulin. Don't forget to connect! improving rag with graph-based reranking. *arXiv preprint arXiv:2405.18414*, 2024.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- Zican Dong, Tianyi Tang, Junyi Li, Wayne Xin Zhao, and Ji-Rong Wen. Bamboo: A comprehensive benchmark for evaluating long text modeling capacities of large language models. *arXiv preprint arXiv:2309.13345*, 2023.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*, 2021.
 URL https://openreview.net/forum?id=XPZIaotutsD.
- 539 Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. Efficient attentions for long document summarization. *arXiv preprint arXiv:2104.02112*, 2021.

554

562

563

564

565 566

567

568

577

- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning. *arXiv preprint arXiv:2112.09118*, 2021.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.
 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. Active retrieval augmented generation. *arXiv preprint arXiv:2305.06983*, 2023.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi
 Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*, 2020.
- Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie
 Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian
 Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina
 Williams. Dynabench: Rethinking benchmarking in nlp, 2021.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
 language models are zero-shot reasoners. *Advances in neural information processing systems*,
 35:22199–22213, 2022.
 - Wojciech Kryściński, Nazneen Rajani, Divyansh Agarwal, Caiming Xiong, and Dragomir Radev. Booksum: A collection of datasets for long-form narrative summarization. arXiv preprint arXiv:2105.08209, 2021.
 - Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. Loogle: Can long-context language models understand long contexts? *arXiv preprint arXiv:2311.04939*, 2023.
- Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. Long-context llms struggle with
 long in-context learning. *arXiv preprint arXiv:2404.02060*, 2024.
- 571
 572
 573
 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004.
- Sheng-Chieh Lin, Akari Asai, Minghan Li, Barlas Oguz, Jimmy Lin, Yashar Mehdad, Wen-tau Yih, and Xilun Chen. How to train your dragon: Diverse augmentation towards generalizable dense retrieval. *arXiv preprint arXiv:2302.07452*, 2023.
- Marina Litvak and Natalia Vanetik. Query-based summarization using mdl principle. In *Proceedings* of the multiling 2017 workshop on summarization and summary evaluation across source types and genres, pp. 22–31, 2017.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and
 Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining
 approach. *CoRR*, abs/1907.11692, 2019. URL http://arxiv.org/abs/1907.11692.
- Zhiyi Ma, Kawin Ethayarajh, Tristan Thrush, Somya Jain, Ledell Wu, Robin Jia, Christopher Potts, Adina Williams, and Douwe Kiela. Dynaboard: An evaluation-as-a-service platform for holistic next-generation benchmarking, 2021.
- Rui Meng, Khushboo Thaker, Lei Zhang, Yue Dong, Xingdi Yuan, Tong Wang, and Daqing He.
 Bringing structure into summaries: a faceted summarization dataset for long scientific documents. arXiv preprint arXiv:2106.00130, 2021.

594 Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and 595 Yoav Shoham. In-context retrieval-augmented language models. Transactions of the Association 596 for Computational Linguistics, 11:1316–1331, 2023. 597 Juan Ramos et al. Using tf-idf to determine word relevance in document queries. In Proceedings of 598 the first instructional conference on machine learning, volume 242, pp. 29–48. Citeseer, 2003. 600 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-601 baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gem-602 ini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530, 2024. 603 604 Stephen Robertson, Hugo Zaragoza, et al. The probabilistic relevance framework: Bm25 and be-605 vond. Foundations and Trends® in Information Retrieval, 3(4):333–389, 2009. 606 607 Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. Ares: An automated evalu-608 ation framework for retrieval-augmented generation systems. arXiv preprint arXiv:2311.09476, 2023. 609 610 Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and 611 Eneko Agirre. Nlp evaluation in trouble: On the need to measure llm data contamination for each 612 benchmark. arXiv preprint arXiv:2310.18018, 2023. 613 Uri Shaham, Maor Ivgi, Avia Efrat, Jonathan Berant, and Omer Levy. Zeroscrolls: A zero-shot 614 benchmark for long text understanding. arXiv preprint arXiv:2305.14196, 2023. 615 616 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya 617 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open 618 models based on gemini research and technology. arXiv preprint arXiv:2403.08295, 2024. 619 karan4d, Teknium, theemozilla, and art huemin. Nous hermes 2 mix-620 dpo. [https://huggingface.co/NousResearch/ tral 8x7b URL 621 Nous-Hermes-2-Mixtral-8x7B-DPO] (https://huggingface.co/ 622 NousResearch/Nous-Hermes-2-Mixtral-8x7B-DPO). 623 624 Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. arXiv 625 preprint arXiv:2212.10509, 2022. 626 627 Alex Wang, Richard Yuanzhe Pang, Angelica Chen, Jason Phang, and Samuel R Bowman. Squality: 628 Building a long-document summarization dataset the hard way. arXiv preprint arXiv:2205.11465, 629 2022. 630 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-631 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language 632 models. arXiv preprint arXiv:2206.07682, 2022a. 633 634 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny 635 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in 636 neural information processing systems, 35:24824–24837, 2022b. 637 Shirley Wu, Shiyu Zhao, Michihiro Yasunaga, Kexin Huang, Kaidi Cao, Qian Huang, Vassilis N. 638 Ioannidis, Karthik Subbian, James Zou, and Jure Leskovec. Stark: Benchmarking llm retrieval 639 on textual and relational knowledge bases, 2024. 640 Peng Xu, Wei Ping, Xianchao Wu, Lawrence McAfee, Chen Zhu, Zihan Liu, Sandeep Subramanian, 641 Evelina Bakhturina, Mohammad Shoeybi, and Bryan Catanzaro. Retrieval meets long context 642 large language models. arXiv preprint arXiv:2310.03025, 2023. 643 644 Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. Benchmarking benchmark leakage in large 645 language models, 2024. 646 Wentao Ye, Jiaqi Hu, Liyao Li, Haobo Wang, Gang Chen, and Junbo Zhao. Data contamination 647 calibration for black-box llms, 2024.

Zichun Yu, Chenyan Xiong, Shi Yu, and Zhiyuan Liu. Augmentation-adapted retriever improves generalization of language models as generic plug-in. arXiv preprint arXiv:2305.17331, 2023. Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. arXiv preprint arXiv:1904.09675, 2019. Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, et al. Extending long context evaluation beyond 100k tokens. arXiv preprint arXiv:2402.13718, 2024. Chenghao Zhu, Nuo Chen, Yufei Gao, and Benyou Wang. Evaluating llms at evaluating temporal generalization, 2024a. Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. Dyval: Dy-namic evaluation of large language models for reasoning tasks. In The Twelfth International Con-ference on Learning Representations, 2024b. URL https://openreview.net/forum? id=gjfOL9z5Xr.

702 A Additional Information for Evalaution

A.1 EVALUATION CRITERIA AND HYPERPARAMETERS

API Access. In this paper, we conduct a comprehensive evaluation over ACADEMICEVAL benchmark using the LLM API provided by together.ai⁶. For each API call, we fix the temperature parameter to 0 (*i.e.*, greedy decoding).

Input Truncation. By default, we use a BERT tokenizer to calculate the number of input tokens for ACADEMICEVAL. However, since the tokenizer of each LLM is usually different, it will cause some inputs to exceed the context length limit of the LLM. Therefore, for the evaluation of each LLM, we additionally download its tokenizer configuration file from the official website at hugging face, which is utilized to ensure correct and accurate truncation of input tokens.

Refinement of LLM Responses. For the TITLE WRITING task, the responses of LLMs are rela-tively short. If the response contains some extra redundant information, it will have a greater impact on the evaluation metric score (although we have given LLM instructions not to generate irrelevant information). Therefore, for the TITLE WRITING task, we additionally refine the LLM responses, for example, removing irrelevant information such as "here is the title". For other tasks, since LLM's responses are relatively long, occasional small amounts of irrelevant information will not have a significant impact on the evaluation, so we do not perform any refinement on LLM's responses in this case.

Details of the implementation of RALM. We use the inputs of ACADEMICEVAL as the external corpus of RALM (such as Target Content and Reference Content introduced in Section F). For text split, we use the RecursiveCharacterTextSplitter from LangChain⁷ and set chunk size and chunk overlap to 512 and 64, respectively. For each retrieval, we recall up to 12 text chunks (limited by the context length of standard LLMs) based on text similarity (semantic similarity based on inner product for dense retrievers or similarity based on word frequency for sparse retrievers).

⁶https://www.together.ai/

⁷https://www.langchain.com/

756 B API COST 757

758	We adopt LLM API provided by together ai ⁸ to conduct experiments in this paper. API costs mainly
759	come from evaluating the test set of ACADEMICEVAL, which are estimated to be around \$300.
760	, , , , , , , , , , , , , , , , , , , ,
761	
762	
763	
764	
765	
766	
767	
768	
769	
770	
771	
772	
773	
774	
775	
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
807	
808	
809	

⁸https://www.together.ai/

LIMITATION AND FUTURE IMPROVEMENT С

ACADEMICEVAL is a live benchmark without label leakage, which leverages co-author papers from a collected co-author graph as few-shot demonstrations to make the context length flexible and scalable. ACADEMICEVAL adopts arXiv as its data source without the need for manual labeling, and the content of the papers on it can naturally serve as high-quality and expert-curated annotations.

However, ACADEMICEVAL still has some limitations:

- **Task Diversity.** ACADEMICEVAL currently has only four academic writing tasks, which limits the task diversity.
- Independent Evaluation of the Paper Section. In ACADEMICEVAL, we independently evaluate the section content extracted from a paper, which may lack a comprehensive evaluation of the paper as a whole.
- Popularity Bias. ACADEMICEVAL first collects a co-author graph from arXiv, which contains a subset of all papers on arXiv. Therefore, the collected papers may have some popularity bias. For example, most of the papers may come from a few active authors, which will cause bias in the evaluation.
- Based on the above limitations, our future improvements will include:
- Introduce More Data Source. The goal of ACADEMICEVAL is to make context length flexible and scalable by using few-shot demonstrations and high-quality labels without manual annotation, so papers on arXiv are a more suitable data source. We will consider adding other websites as data sources in the future, such as some question-answering websites (Stack Overflow or Reddit, etc.). In this case, we can use the best answers as high-quality labels. By modeling the citation relationship between posts into a graph, we can also obtain few-shot demonstrations to enrich the context length.
 - K-fold Cross-validation. We can use k-fold cross-validation for a paper, that is, leaving a section (or fold) as the label, the remaining sections as inputs, and finally calculating the average of all leave-one-out evaluation scores.
 - Eliminate Popularity Bias. We will perform probabilistic sampling on papers when collecting the co-author graph and give a lower sampling probability to active authors to alleviate the impact of popularity bias.

864 D SOCIAL IMPACT

The proposed benchmark ACADEMICEVAL will promote the academic community's exploration of using LLMs to automate academic writing tasks. Here are some key points highlighting its significance:

- Efficiency and Productivity. LLMs can drastically reduce the time and effort required for various academic writing tasks. These tasks include drafting papers, writing literature reviews, summarizing research articles, and generating bibliographies. By automating these processes, researchers can focus more on high-level thinking, experimentation, and analysis.
- Enhanced Writing Quality. LLMs have the ability to produce coherent and grammatically correct text, which can improve the overall quality of academic writing. They can assist in refining arguments, improving clarity, and ensuring consistency in style and tone, which is particularly useful for non-native English speakers.
- **Support for Multidisciplinary Research.** Given their training on diverse topics, LLMs can assist researchers in exploring interdisciplinary approaches by providing information and generating content across various fields of study. This can foster innovation and encourage collaboration between different academic disciplines.

Е DETAILS OF LIVE LEADERBOARD

To enhance usability, we aim to create a live leaderboard on Hugging Face to help users easily utilize our benchmarks and compare various models. Our leaderboard will provide the following functionalities:

- 1. Live updates and time selection: We will update the dataset periodically to ensure it includes the latest papers from arXiv, which the LLMs have never seen before, to prevent label leakage. Users can choose the version they wish to use.
- 2. Different abstraction tasks: We will create separate leaderboards for each task. Users are welcome to run their models on one or more tasks and report their results.
 - 3. Ease of use: We will provide detailed and standardized instructions so that users can easily run the pipeline with their models and obtain the results.

972 F LLM PROMPTS

974 975

976

977

978

979

In this section, we present the LLM prompts used in the experiments, including TITLE WRIT-ING, ABSTRACT WRITING, INTRODUCTION WRITING, and RELATED WORK WRITING. For each academic writing task, we provide prompts for standard LLMs, long-context LLMs, and RALM (RALM additionally includes the retrieval query).

F.1 LLM PROMPTS FOR TITLE WRITING

986 987

989 990

991

992

993

994

995

996

997

998

999 1000

1002

1004

Prompt for Standard and Long-context LLMs on TITLE-10K

Please read the following Target Content carefully and summarize the Target Content as required.

Target Content: {CONTENT}

Target Content Abstract: {ABSTRACT}

Please craft a title highly summarizing the main theme from the above provided Target Content. The title should be of appropriate length (strictly limited to about 10 words). The title should also include and highlight the core and most critical theme of the Target Content, ignoring minor and redundant information. Please ensure that the title captures the essence of the Target Content in a clear and concise manner. Please output the title directly without including other redundant or irrelevant text.

Prompt for Standard and Long-context LLMs on TITLE-30K and TITLE-31K-G

1008 Please read the following Reference Content and Output carefully and summarize the Target 1009 Content as required. 1010 ### Reference Content 0: {CONTENT_0} 1011 ### Reference Abstract 0: {ABSTRACT_0} 1012 ### Reference Output 0: {OUTPUT_0} 1013 1014 ### Target Content: {CONTENT} 1015 ### Target Content Abstract: {ABSTRACT} 1016 Please craft a title highly summarizing the main theme from the above provided Target Content. The Reference Content and Output provide some demonstrations, which may also 1017 contain some information that is potentially related to the Target Content. You can refer to the input and output text forms of the Reference Content and Output to assist in summarizing the Target Content and try to explore and use the information that is potentially related to the Target Content contained in the Reference Content and Output. The title should be 1021 of appropriate length (strictly limited to about 10 words). The title should also include and highlight the core and most critical theme of the Target Content, ignoring minor and redun-1023 dant information. Please ensure that the title captures the essence of the Target Content in a 1024 clear and concise manner. Please output the title directly without including other redundant 1025 or irrelevant text.

Prom	pt for RALM on TITLE-10K, TITLE-30K, and TITLE-31K-G
Pleas	e read the following Target Content carefully and summarize the Target Content as
requi	red.
ן' ### ר ### יו	arget Content 0: {CONTENT_0}
### 1	
### T	`arget Content Abstract: {ABSTRACT}
Pleas	e craft a title highly summarizing the main theme from all the above provided Target
Conte The t	it is should also include and highlight the core and most critical theme of the Target
Conte	ents, ignoring minor and redundant information. Please ensure that the title captures
the es	ssence of the Target Contents in a clear and concise manner. Please output the title
direct	ly without including other redundant or irrelevant text.
Potri	aval Quary for PAIM on TITLE 10K TITLE 20K and TITLE 21K C
Ketti	eval Query for KALWI on TITLE-TOK, TITLE-SOK, and TITLE-STK-O
Pleas the te	e craft a title highly summarizing the main theme of the provided text. The abstract of xt is: {ABSTRACT}
2 LL	M PROMPTS FOR ABSTRACT WRITING
Prom	pt for Standard and Long-context LLMs on ABS-9K
Pleas	e read the following Target Content carefully and summarize the Target Content as
requi	red.
###]	arget Content: {CONTENT}
### 1 Pleas	arget Content 11tle: { 111LE } e craft an abstract summarizing the key points from the above provided Target Content
The a	bstract should be of appropriate length (around 200 words) and include the main theme.
signif	icant findings or arguments, and conclusions of the Target Content. Please ensure that
the a	ostract captures the essence of the Target Content in a clear, coherent, and succinct
mann text	er. Please output the abstract directly without including other redundant or irrelevant
Prom	nt for Standard and Long-context LLMs on ABS-28K and ABS-29K-G
110111	r for building and Long context ELIVIS on The 20th and The 27th G
Pleas	e read the following Reference Content and Output carefully and summarize the Target
Conte	int as required.
### F	Reference Title 0: {TITLE 0}
### F	Reference Output 0: {OUTPUT_0}
###]	arget Content: {CONTENT}
###] DI	arget Content Title: {TITLE}
Pleas	e craft an abstract summarizing the key points from the above provided Target Content.
some	information that is potentially related to the Target Content. You can refer to the input
and c	butput text forms of the Reference Content and Output to assist in summarizing the
Targe	t Content and try to explore and use the information that is potentially related to the
Targe	t Content contained in the Reference Content and Output. The abstract should be of
appro	priate length (around 200 words) and include the main theme, significant findings or
argun	nents, and conclusions of the Target Content. Please ensure that the abstract captures
the el	sence of the larget Content in a clear, concrent, and succinct manner. Please output
un di	shaet anothy without motuling other requilitant of inclosent text.

1080	
1081	Prompt for RALM on ABS-10K, ABS-30K, and ABS-31K-G
1082	Please read the following Target Content carefully and summarize the Target Content as re-
1083	anired
1084	### Target Content 0: {CONTENT_0}
1085	### Target Content 1: {CONTENT_1}
1086	
1087	### Target Content Title: {TITLE}
1088	Please craft an abstract summarizing the key points from all the above provided Target Con-
1089	tents. The abstract should be of appropriate length (around 200 words) and include the main
1090	theme, significant findings or arguments, and conclusions of the Target Contents. Please
1091	ensure that the abstract captures the essence of the Target Contents in a clear, conferent, and succinet manner. Please output the abstract directly without including other redundant or
1092	irrelevant text
1093	
1094	
1095	
1096	
1097	
1098	
1099	
1100	
1101	
1102	Retrieval Query for RALM on ABS-10K, ABS-30K, and ABS-31K-G
1103	Please craft an abstract summarizing the key points of the provided text. The title of the text
1104	i lease cruit an abstract summarizing the key points of the provided text. The the of the text

F.3 LLM PROMPTS FOR INTRODUCTION WRITING

1104

1116

is: {TITLE}

Prompt for Standard and Long-context LLMs on INTRO-8K

1117 Please read the following Target Content carefully and summarize the Target Content as re-1118 quired. 1119 ### Target Content: {CONTENT} 1120 ### Target Content Title: {TITLE} 1121 ### Target Content Abstract: {ABSTRACT} Please craft an introduction summarizing the key points from the above provided Target 1122 Content. The introduction should be of appropriate length (about 1000 to 1500 words). The 1123 introduction should first describe the topic or main theme of the Target Content, then provide 1124 relevant background knowledge, and summarize the existing relevant research on this topic 1125 from the Target Content, point out their advantages and disadvantages, and highly summa-1126 rize the specific research problem and problem statement targeted by the Target Content. 1127 Next, describe in detail the core approach or insights proposed by the Target Content on this 1128 topic and include any necessary experimental results. Then, use about 3 short paragraphs 1129 (each paragraph is about 50 words) to highly summarize the approach or insights proposed 1130 in the Target Content, as well as the experimental results. Finally, briefly give an overview 1131 of the Target Content's structure. Please ensure that the introduction captures the essence of 1132 the Target Content in a clear, coherent, and succinct manner. Please output the introduction 1133 directly without including other redundant or irrelevant text.

1134	Prompt for Standard and Long-context LLMs on INTRO-28K and INTRO-28K-G
1135	
1130	Please read the following Reference Content and Output carefully and summarize the Target
1107	Content as required.
1120	### Reference Content 0: {CONTENT_0}
1139	### Reference Abstract 0: $\{ABSTRACT 0\}$
1140	### Reference Output 0: {OUTPUT 0}
1142	
1143	### Target Content: {CONTENT}
1144	### Target Content Title: {TITLE}
1145	### Target Content Abstract: {ABSTRACT}
1146	Please craft an introduction summarizing the key points from the above provided Target
1147	Content. The Reference Content and Output provide some demonstrations, which may also
1148	contain some information that is potentially related to the larget Content. You can refer to
1140	the input and output text forms of the Reference Content and Output to assist in summarizing
1150	the Target Content and if y to explore and use the information that is potentially related to the Target Content contained in the Reference Content and Output. The introduction should
1151	be of appropriate length (about 1000 to 1500 words). The introduction should first describe
1152	the topic or main theme of the Target Content, then provide relevant background knowledge,
1153	and summarize the existing relevant research on this topic from the Target Content, point
1154	out their advantages and disadvantages, and highly summarize the specific research problem
1155	and problem statement targeted by the Target Content. Next, describe in detail the core
1156	approach or insights proposed by the Target Content on this topic and include any necessary
1157	experimental results. Then, use about 3 short paragraphs (each paragraph is about 50 words)
1158	to highly summarize the approach or insights proposed in the larget Content, as well as
1159	Please ensure that the introduction cantures the essence of the Target Content in a clear
1160	coherent, and succinct manner. Please output the introduction directly without including
1161	other redundant or irrelevant text.
1162	
1163	
1164	
1165	
1166	
1167	
1168	Prompt for KALM on INTRO-8K, INTRO-28K, and INTRO-28K-G
1169	Please read the following Target Content carefully and summarize the Target Content as re-
1170	quired.
1171	### Target Content 0: {CONTENT_0}
1172	### Target Content 1: {CONTENT_1}
1173	
1174	### Target Content Title: {TITLE}
1175	### Target Content Abstract: {ABSTRACT}
1176	Please craft an introduction summarizing the key points from all the above provided Target
1177	Contents. The introduction should be of appropriate length (about 1000 to 1500 words). The
1178	nitroduction should first describe the topic of main theme of the faiget Contents, then pro-
1179	from the Target Contents, point out their advantages and disadvantages, and highly summa-
1180	rize the specific research problem and problem statement targeted by the Target Contents
1181	Next, describe in detail the core approach or insights proposed by the Target Contents on this
1182	topic and include any necessary experimental results. Then, use about 3 short paragraphs
1183	(each paragraph is about 50 words) to highly summarize the approach or insights proposed
1184	in the Target Contents, as well as the experimental results. Finally, briefly give an overview
1185	of the Target Contents' structure. Please ensure that the introduction captures the essence of
1186	the Target Contents in a clear, coherent, and succinct manner. Please output the introduction
1187	directly without including other redundant or irrelevant text.

	Keuleval Query for KALW on INTRO-6K, INTRO-26K, and INTRO-26K-O
	Please craft an introduction summarizing the main theme of the provided text (including background knowledge, advantages and disadvantages of existing research and challenges, the proposed approach, experimental results, etc.). The title of the text is {TITLE}. The abstract of the text is {ABSTRACT}.
7	4 II M DROMPTS FOR DELATED WORK WRITING
	4 LLM FROMPTS FOR RELATED WORK WRITING
	Prompt for Standard and Long-context LLMs on RELATED-34K
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. #### Target Citation 0:
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0}
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: [CONTENT]
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE}
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Abstract: {ABSTRACT} Cited to The of Content target the theorem IT the charget Cited in the first of the time for the time
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Title: {ABSTRACT} Given the Target Content and its Abstract and Title, along with its Target Citations (including Target Citation Title and Abstract), please craft a related work summarizing the key points
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Title: {ABSTRACT} Given the Target Content and its Abstract and Title, along with its Target Citations (including Target Citation Title and Abstract), please craft a related work summarizing the key points from the above provided Target Citations. There is no specific length requirement or limit
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Title: {ABSTRACT} Given the Target Content and its Abstract and Title, along with its Target Citations (including Target Citation Title and Abstract), please craft a related work summarizing the key points from the above provided Target Citations. There is no specific length requirement or limit for the entire related work (it is best to keep it around 500 to 1000 words), but each Target Citation Line and Abstract has the label of the la
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Title: {TITLE} ### Target Content Abstract: {ABSTRACT} Given the Target Content and its Abstract and Title, along with its Target Citations (including Target Citation Title and Abstract), please craft a related work summarizing the key points from the above provided Target Citations. There is no specific length requirement or limit for the entire related work (it is best to keep it around 500 to 1000 words), but each Target Citation that appears in the related work needs to be highly summarized in extremely concise and short sentances. You can refer to the topic or main theme described by the Target
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Title: {TITLE} ### Target Content Abstract: {ABSTRACT} Given the Target Content and its Abstract and Title, along with its Target Citations (including Target Citation Title and Abstract), please craft a related work summarizing the key points from the above provided Target Citations. There is no specific length requirement or limit for the entire related work (it is best to keep it around 500 to 1000 words), but each Target Citation that appears in the related work needs to be highly summarized in extremely concised and short sentences. You can refer to the topic or main theme described by the Target Content and its Abstract and Title to filter irrelevant information in the Target Citations and the target Citations and the target citations and Please to the topic or main theme described by the Target Content and its Abstract and Title to filter irrelevant information in the Target Citations and Please to the topic or main theme described by the Target Content and its Abstract and Title to filter irrelevant information in the Target Citations and Please to the topic or main theme described by the Target Please to the topic or main theme described by the Target Please to the topic or main theme described by the Target Please to the topic or main theme described by the Target Please to the topic or main theme described by the Target Please to the topic or main theme topic the topic or main theme topic the topic or the to
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Title: {TITLE} ### Target Content Abstract: {ABSTRACT} Given the Target Content and its Abstract and Title, along with its Target Citations (including Target Citation Title and Abstract), please craft a related work summarizing the key points from the above provided Target Citations. There is no specific length requirement or limit for the entire related work (it is best to keep it around 500 to 1000 words), but each Target Citation that appears in the related work needs to be highly summarized in extremely concises and short sentences. You can refer to the topic or main theme described by the Target Content and its Abstract and Title to filter irrelevant information in the Target Citations and leverage relevant information. Furthermore, you can categorize the relevant Target Citations
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Title: {TITLE} ### Target Content Abstract: {ABSTRACT} Given the Target Content and its Abstract and Title, along with its Target Citations (including Target Citation Title and Abstract), please craft a related work summarizing the key points from the above provided Target Citations. There is no specific length requirement or limit for the entire related work (it is best to keep it around 500 to 1000 words), but each Target Citation that appears in the related work needs to be highly summarized in extremely concises and short sentences. You can refer to the topic or main theme described by the Target Content and its Abstract and Title to filter irrelevant information in the Target Citations and leverage relevant information. Furthermore, you can categorize the relevant Target Citations briefly summarize the advantages and disadvantages of each categorization, and explain the
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C_TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Title: {TITLE} ### Target Content Abstract: {ABSTRACT} Given the Target Content and its Abstract and Title, along with its Target Citations (including Target Citation Title and Abstract), please craft a related work summarizing the key points from the above provided Target Citations. There is no specific length requirement or limit for the entire related work (it is best to keep it around 500 to 1000 words), but each Target Citation that appears in the related work needs to be highly summarized in extremely concises and short sentences. You can refer to the topic or main theme described by the Target Content and its Abstract and Title to filter irrelevant information in the Target Citations and leverage relevant information. Furthermore, you can categorize the relevant Target Citations and leverage relevant information. Furthermore, you can categorize the relevant Target Citations briefly summarize the advantages and disadvantages of each categorization, and explain the advantages of the approach proposed in the Target Content. Please ensure that the related work cantures all the relevant key noints of the Target Citations in a clear coherent and advantages of the approach proposed in the Target Citations in a clear coherent.
	Prompt for Standard and Long-context LLMs on RELATED-34K Please read the following Target Content and Target Citations carefully and summarize the Target Citations according to the topic of the Target Content as required. ### Target Citation 0: Target Citation Title: {C.TITLE_0} Target Citation Abstract: {C_ABSTRACT_0} ### Target Content: {CONTENT} ### Target Content Title: {TITLE} ### Target Content Title: {TITLE} ### Target Content Abstract: {ABSTRACT} Given the Target Content and its Abstract and Title, along with its Target Citations (including Target Citation Title and Abstract), please craft a related work summarizing the key points from the above provided Target Citations. There is no specific length requirement or limit for the entire related work (it is best to keep it around 500 to 1000 words), but each Target Citation that appears in the related work needs to be highly summarized in extremely concise and short sentences. You can refer to the topic or main them described by the Target Content and its Abstract and Title to filter irrelevant information in the Target Citations, briefly summarize the advantages and disadvantages of each categorization, and explain the advantages of the approach proposed in the Target Content. Please ensure that the related work captures all the relevant key points of the Target Citations in a clear, coherent, and succinct manner. Please output the related work directly without including other redundant

1242	Prompt for Standard and Long-context LLMs on RELATED-53K and RELATED-53K-G
1243	
1244	Please read the following Reference Content and Output carefully and summarize the Target
1245	Citations according to the topic of the Target Content as required.
1246	### Reference Content 0: {CONTENT_0}
1247	$### \text{Reference Title 0: } \{TTTLE_0\}$
1248	### Reference Abstract 0: {ABSTRACT_0}
1249	$\#\#\#$ Reference Output 0: $\{OO IPO I_0\}$
1250	 ### Target Citation 0:
1251	Target Citation Title: $\{C \text{ TITLE } 0\}$
1252	Target Citation Abstract: {C ABSTRACT 0}
1253	
1254	### Target Content: {CONTENT}
1255	### Target Content Title: {TITLE}
1256	### Target Content Abstract: {ABSTRACT}
1257	Given the Target Content and its Abstract and Title, along with its Target Citations (including
1258	Target Citation Title and Abstract), please craft a related work summarizing the key points
1259	from the above provided Target Citations. The Reference Content and Output provide some
1260	demonstrations, which may also contain some information that is potentially related to the
1261	Target Content. You can refer to the input and output text forms of the Reference Content
1262	and Output to assist in summarizing the Target Citations and try to explore and use the in-
1263	formation (e.g., related citations missing from the Target Citations) that is potentially related
1264	to the Target Content contained in the Reference Content and Output. There is no specific
1265	length requirement or limit for the entire related work (it is best to keep it around 500 to 1000
1266	words), but each Target Chatton that appears in the related work needs to be highly summa-
1267	described by the Target Content and its Abstract and Title to filter irrelevant information
1268	in the Target Citations and leverage relevant information. Furthermore, you can categorize
1269	the relevant Target Citations briefly summarize the advantages and disadvantages of each
1270	categorization, and explain the advantages of the approach proposed in the Target Content.
1271	Please ensure that the related work captures all the relevant key points of the Target Citations
1272	in a clear, coherent, and succinct manner. Please output the related work directly without
1273	including other redundant or irrelevant text.
1274	
1975	
1275	
1277	
1079	Prompt for RALM on RELATED-34K, RELATED-53K, and RELATED-53K-G
1270	
12/9	Please read the following Target Content and Target Citations carefully and summarize the
1200	Target Citations according to the topic of the Target Content as required.
1201	### Target Content 0: {CONTENT_0}
1282	### Target Content 1: {CONTENT_1}
1283	
1284	### Target Content Title: {TITLE}
1285	### Target Content Abstract: {ABSTRACT}
1286	I Given the Target Content Abstract and Litle place craft a related work summarizing the key

Given the Target Content Abstract and Title, please craft a related work summarizing the key points from all the above provided Target Contents. There is no specific length requirement 1287 or limit for the entire related work (it is best to keep it around 500 to 1000 words), but each 1288 Target Content that appears in the related work needs to be highly summarized in extremely 1289 concise and short sentences. You can refer to the topic or main theme described by the Target 1290 Content Abstract and Title to filter irrelevant information in the Target Contents and leverage 1291 relevant information. Furthermore, you can categorize the relevant Target Contents and 1292 briefly summarize the advantages and disadvantages of each categorization. Please ensure 1293 that the related work captures all the relevant key points of the Target Contents in a clear, 1294 coherent, and succinct manner. Please output the related work directly without including 1295 other redundant or irrelevant text.

1296	Retrieval Query for RALM on RELATED-34K, RELATED-53K, and RELATED-53K-G
1297	
1298	Please craft a related work summarizing all the relevant key points of the provided text. The
1299	title of the text is {TITLE}. The abstract of the text is {ABSTRACT}.
1300	
1301	
1302	
1303	
1304	
1305	
1306	
1307	
1308	
1309	
1011	
1010	
1212	
1217	
1315	
1316	
1317	
1317	
1310	
1320	
1321	
1322	
1323	
1324	
1325	
1326	
1327	
1328	
1329	
1330	
1331	
1332	
1333	
1334	
1335	
1336	
337	
1338	
1339	
1340	
1341	
1342	
1343	
1344	
1345	
1346	
1347	
1348	
1349	

1350 G EXAMPLES

551	
352	
53	
54	In this section, we present some generation examples w.r.t. TITLE WRITING, ABSTRACT WRITING,
5	and INTRODUCTION WRITING tasks.
	G.1 TITLE WRITING
1	Comparison between Generated Title and Ground Truth on TITLE WRITING Task
	r
	Ground Truth:
	When to Trust LLMs: Aligning Confidence with Response Quality
	• Qwen 1.5 Chat (72B)
	Commented Title (0, shet):
	Generated Little (U-snot):
	Language Models
	Generated Title (2-shot demonstrations from random selection):
	CONOORD: Enhancing LLM Confidence Alignment for Improved Trustworthiness and
	Retrieval
	Generated Title (2-shot demonstrations from co-author papers):
	CONQORD: Enhancing Confidence Calibration for Language Models via Reinforcement
	Learning
	· Mirtual 8x7D Instruct (167D)
	• \mathbf{W} \mathbf{X} \mathbf{U}
	Generated Title (0-shot):
	CONOORD: Aligning Confidence with Quality in Large Language Models
	Generated Title (2-shot demonstrations from random selection):
	CONQORD: Reinforcement Learning for Confidence-Quality Alignment in LLMs
	Generated Title (2-shot demonstrations from co-author papers):
	CONCORD: Angning Confidence with Response Quality in LLWIS
	• Nous Hermes 2 - Mixtral 8x7B-DPO (46 7B)
	$= 110005 11010105 \Delta = 1111101010A / D^2 D I O (100/D)$
ļ	Generated Title (0-shot):
	Aligning Confidence with Response Quality in Large Language Models using Reinforce-
	ment Learning with Dual-Component Reward Function: CONQORD Approach
	Generated Title (2-shot demonstrations from random selection):
	Aligning Confidence with Response Quality in Large Language Models: Introducing
	CONQORD
	Concreted Title (2 shot domonstrations from as suffice reasons).
	Confidence Alignment in Large Language Models via Deinforcement Learning: A CON
	OORD Approach for Trust
	Yord Approach for fluor

	Comparison between Generated Title and Ground Truth on TITLE WRITING Task
	Cround Truth
	Ground fruit: Distillation for Multilingual Information Patriaval
	Distination for Mutuningual mormation Retrieval
	• Owen 1.5 Chat (72B)
	Generated Title (0-shot):
	Multilingual Translate-Distill: Enhancing Cross-Language Information Retrieval Models
	Generated Title (2-shot demonstrations from random selection):
Wullingual Translate-Distill for Enhanced Cross-Language Information Retriev	Multilingual Translate-Distill for Enhanced Cross-Language Information Retrieval
Generated Title (2-shot demonstrations from co-author papers): Multilingual Translate-Distill for Effective Multilingual Information Retrieval	Generated Title (2-shot demonstrations from co-author naners):
	Multilingual Translate-Distill for Effective Multilingual Information Retrieval
	Mataningaa maistae Distin for Enective mataningaa mornadon teateva
	• Mixtral-8x7B Instruct (46.7B)
	Generated Title (0-shot):
Multilingual Translate-Distill for Cross-Language Information Retrieval	Multilingual Translate-Distill for Cross-Language Information Retrieval
Generated Title (2-shot demonstrations from random selection):	Generated Title (2-shot demonstrations from random selection):
	Multilingual Translate-Distill for Cross-Language Information Retrieval
Concreted Title (2 shot domonstrations from as without more as)	Concrated Title (2-shot demonstrations from co-author papars).
	Multilingual Translate-Distill for MLR. Extending Translate-Distill for Robust Language
Mixing	Mixing Manshale Distin for MERC. Exceloring Hanshale-Distin for Robust Language
	0
	Nous Hermes 2 - Mixtral 8x7B-DPO (46.7B)
	Generated Title (0-shot):
	Multilingual Translate-Distill: Enhancing Multilingual Information Retrieval with Transla-
	tion and Distillation
	Concreted Title (2-shot demonstrations from random selection).
	Multilingual Translate-Distill: Efficient Training for Multilingual Information Retrieval
	Translate District Different Training for Wutchinguar information Refleval
	Generated Title (2-shot demonstrations from co-author papers):
Multilingual Translate-Distill: Training Dual-Encoder Models for Multilingual Informa	Multilingual Translate-Distill: Training Dual-Encoder Models for Multilingual Information
	Retrieval

1458 G.2 Abstract Writing

Comparison between Generated Abstract and Ground Truth on ABSTRACT WRITING Task

Ground Truth:

In recent years, neural ranking models (NRMs) have been shown to substantially outperform their lexical counterparts in text retrieval. In traditional search pipelines, a combination of features leads to well-defined behavior. However, as neural approaches become increasingly prevalent as the final scoring component of engines or as standalone systems, their robustness to malicious text and, more generally, semantic perturbation needs to be better understood. We posit that the transformer attention mechanism can induce exploitable defects through positional bias in search models, leading to an attack that could generalize beyond a single query or topic. We demonstrate such defects by showing that non-relevant text-such as promotional content-can be easily injected into a document without adversely affecting its position in search results. Unlike previous gradient-based attacks, we demonstrate these biases in a query-agnostic fashion. In doing so, without the knowledge of topicality, we can still reduce the negative effects of non-relevant content injection by controlling injection position. Our experiments are conducted with simulated on-topic promotional text automatically generated by prompting LLMs with topical context from target documents. We find that contextualization of a non-relevant text further reduces negative effects whilst likely circumventing existing content filtering mechanisms. In contrast, lexical models are found to be more resilient to such content injection attacks. We then investigate a simple yet effective compensation for the weaknesses of the NRMs in search, validating our hypotheses regarding transformer bias.

• Qwen 1.5 Chat (72B)

Generated Abstract (2-shot demonstrations from co-author papers):

The study investigates the vulnerability of Neural Ranking Models (NRMs) to adversarial attacks through the injection of promotional content into documents. The authors propose a framework for generating contextualised text using Large Language Models (LLMs) and examine the impact of position and context on the relevance estimation of injected text. They find that NRMs, particularly those based on BERT and T5 architectures, are sensitive to the position of injected text, with content placed after salient sentences having a lesser negative effect on relevance. The study also introduces a mitigation strategy using a classifier to detect promotional content, demonstrating its effectiveness in reducing the negative impact on retrieval performance. The findings highlight the potential for exploiting positional bias in NRMs and the need for robustness against such attacks in semantic search.

• Nous Hermes 2 - Mixtral 8x7B-DPO (46.7B)

Generated Abstract (2-shot demonstrations from co-author papers):

This study investigates the exploitation of positional bias in neural ranking models (NRMs) for query-agnostic generative content injection in search. The authors propose the concept of attention bleed-through, where positive attention scores from relevant spans propagate to non-relevant spans, affecting relevance estimation. They find that injecting non-relevant text close to document sentences with high similarity scores with the query can largely affect the rank of documents. The study also demonstrates that controlling the position and contextualization of promotional text can mitigate its negative impact on retrieval performance. A model-agnostic mitigation strategy is proposed, which improves nDCG@10 significantly under a classic evaluation setting by reducing the effect of contextualization. The findings highlight the importance of addressing positional bias in NRMs and the potential risks of generative content injection in semantic search.

Comparison between Generated Abstract and Ground Truth on ABSTRACT WRITING Task

Ground Truth:

1512

1513 1514

1515

1516

1517

1518

1519

1520

1521

1522

1525

1527

1528 1529

1530

1531

1532

1533

1534

1535

1536

1537

1538

1540

1541

1542

1543

1544 1545

1546 1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

Developing autonomous agents for mobile devices can significantly enhance user interactions by offering increased efficiency and accessibility. However, despite the growing interest in mobile device control agents, the absence of a commonly adopted benchmark makes it challenging to quantify scientific progress in this area. In this work, we introduce B-MoCA, a novel benchmark designed specifically for evaluating mobile device control agents. To create a realistic benchmark, we develop B-MoCA based on the Android operating system and define 60 common daily tasks. Importantly, we incorporate a randomization feature that changes various aspects of mobile devices, including user interface layouts and language settings, to assess generalization performance. We benchmark diverse agents, including agents employing large language models (LLMs) or multi-modalLLMs as well as agents trained from scratch using human expert demonstrations. While these agents demonstrate proficiency in executing straightforward tasks, their poor performance on complex tasks highlights significant opportunities for future research to enhance their effectiveness. Our source code is publicly available at https://b-moca.github.io.

• Qwen 1.5 Chat (72B)

Generated Abstract (2-shot demonstrations from co-author papers):

The Benchmark for Mobile Device Control Agents (B-MoCA) is a novel evaluation platform designed to assess the performance of agents in controlling mobile devices across diverse configurations. The benchmark features 60 practical tasks grounded in real-life scenarios and supports randomization of device attributes like icon placements, wallpapers, languages, and device types to test generalization ability. The study compares three types of agents: LLM (Large Language Model) agents, MLLM (Multimodal Large Language Model) agents, and Vision-Language-UI (VLUI) agents. Results show that while LLM and MLLM agents perform well in simple tasks, they struggle with complex tasks requiring sequential decision-making. VLUI agents, trained from scratch, demonstrate better performance in completing tasks but still face challenges with severe visual changes in unseen device configurations. The study highlights the importance of diverse training data and the need for more efficient methods to leverage multi-modal inputs. Overall, B-MoCA provides a comprehensive framework for evaluating and advancing mobile device control agents in real-world settings.

• Nous Hermes 2 - Mixtral 8x7B-DPO (46.7B)

Generated Abstract (2-shot demonstrations from co-author papers):

This paper introduces B-MoCA, a benchmark designed to evaluate the performance of mobile device control agents on diverse device configurations in executing common daily tasks. The benchmark includes 60 basic tasks grounded in realistic scenarios and a randomization feature that changes icon placements and sizes, wallpapers, languages, and device types. Baseline agents include LLM agents, MLLM agents, and Vision-Language-UI (VLUI) agents. Experiments reveal that LLM agents and MLLM agents show high robustness across diverse device configurations while they fall short on multiple sequential decision-making. Agents trained with BC successfully mimic expert behaviors but lack generalization ability in test environments with unseen device configurations. The study highlights the limitations of existing methods in mobile device control and calls for future research.

- 559
- 1560
- 1561
- 1563
- 1564
- 1565

1566 G.3 INTRODUCTION WRITING

1579

1580

1581

1584 1585

Comparison between Generated Introduction and Ground Truth on INTRODUCTION WRIT-ING Task

Ground Truth:

The original text is too long to be included here. Please see 2404.17196v1 via arXiv for more details.

• Qwen 1.5 Chat (72B)

Generated Introduction (2-shot demonstrations from co-author papers):

1586 Retrieval poisoning is a newly identified threat to large language model (LLM)-powered 1587 applications, particularly those utilizing the retrieval augmented generation (RAG) tech-1588 nique. This threat arises from the vulnerability in LLM application frameworks, which allow attackers to subtly manipulate the content of documents in a way that appears benign to humans but can guide the LLMs to generate incorrect or malicious responses. 1590 The attack is executed by injecting invisible attack sequences into the documents, which 1591 remain undetected during the RAG process, leading to compromised application responses. 1592 The prevalence of LLM application frameworks, such as LangChain, has enabled the development of numerous LLM-powered applications. However, these frameworks often 1594 overlook the security risks associated with external content. Retrieval poisoning exploits this vulnerability by crafting malicious documents that, despite containing accurate information, 1596 can prompt the LLMs to produce misleading responses when used as reference sources during RAG. The success rate of such attacks is alarming, with preliminary experiments 1598 demonstrating an 88.33% success rate in misleading LLMs and a 66.67% success rate in real-world applications. To execute retrieval poisoning, attackers analyze the LLM application framework, identifying exploitable components like document parsers, text splitters, and prompt templates. They then generate an attack sequence that is invisible to human perception but influences the LLM's response generation. The crafted malicious document is designed to maintain its original appearance while containing the hidden attack sequence, which is strategically placed to remain undivided by text splitters. The proposed 1604 approach involves a two-step process: analyzing the framework to identify invisible injec-1605 tion points and generating the attack sequence. Attackers leverage the framework's prompt 1606 templates to create augmented requests that guide the LLMs towards generating targeted malicious responses. The effectiveness of the attack is demonstrated through experiments on different LLMs, augmented requests, and real-world applications, showcasing the 1609 potential impact of retrieval poisoning. In summary, retrieval poisoning is a stealthy attack 1610 method that exploits the vulnerabilities in LLM application frameworks, particularly in the RAG process. By injecting invisible attack sequences into documents, attackers can 1611 manipulate LLMs into generating incorrect responses, posing a significant security concern 1612 for LLM-powered applications. The proposed approach highlights the need for better 1613 understanding of LLM application frameworks and the development of more robust defense 1614 mechanisms to mitigate such threats. The paper structure delves into the problem, presents 1615 the attack methodology, demonstrates its effectiveness through experiments, and discusses 1616 potential defense strategies, ultimately contributing to the ongoing research on securing 1617 LLM-powered applications. 1618

1620	Comparison between Generated Introduction and Ground Truth on INTRODUCTION WRIT-
1621	ING Task
1622	
1623	Ground Truth:
1624	The original text is too long to be included here. Please see 2404.18465v2 via arXiv for
1625	more details.
1626	
1627	• Qwen 1.5 Chat (72B)
1628	Concreted Introduction () shot domenstrations from as outhor nanous).
1629	The M3oE framework is a pioneering solution to the challenge of multi domain multi task
1630	recommendation addressing the intricate dependencies among multiple domains and tasks
1631	in a disentangled and adaptive manner. Recognizing the need for a more comprehensive
1632	approach in practical recommendation systems, M3oE integrates multi-domain information
1633	and leverages knowledge transfer across domains and tasks, optimizing for multiple
1634	objectives simultaneously. The framework consists of three mixture-of-experts modules,
1635	each dedicated to capturing common, domain-specific, and task-specific user preferences.
1636	A two-level fusion mechanism is employed to effectively extract and integrate features
1637	across diverse domains and tasks, ensuring a balanced representation of user preferences.
1638	Existing multi-domain and multi-task recommendation methods have shown promise in
1639	enhancing recommendation accuracy by exploiting shared patterns. However, they often
1640	fall short when faced with the complexity of simultaneously addressing multiple domains
1641	to optimize the model structure dynamically. This adaptability is crucial for real world
1642	scenarios where data characteristics and task requirements can vary significantly. The core
1643	of M3oE lies in its three expert modules: shared domain and task experts. The shared
1644	expert module learns common patterns across domains and tasks, while the domain expert
1645	module focuses on domain-specific characteristics. The task expert module, on the other
1646	hand, models task-specific user preferences. The two-level fusion mechanism balances
1647	the contribution of these modules, allowing for a flexible and effective representation
1648	of user preferences. Experimental results on benchmark datasets demonstrate M3oE's
1649	superior performance compared to state-of-the-art baselines, highlighting its ability to
1650	address the multi-domain multi-task seesaw problem. In summary, M3oE is a novel and
1651	adaptive multi-domain multi-task recommendation framework that integrates multi-domain
1652	chieve and the second s
1653	AutoMI driven fusion weights enable it to effectively capture user preferences in complex
1654	recommendation scenarios. The proposed approach not only surpasses existing methods in
1655	performance but also offers a more generalizable solution to the multi-domain multi-task
1656	recommendation challenge. The structure of the paper delves into the problem definition.
1657	framework overview, detailed module descriptions, experimental setup, and results, provid-
1658	ing a comprehensive understanding of M3oE's design and effectiveness.
1659	
1660	