

XL²Bench: A Benchmark for Extremely Long Context Understanding with Long-range Dependencies

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have demonstrated remarkable performance across diverse tasks but are constrained by their small context window sizes. Various efforts have been proposed to enhance the capability of LLMs to process and comprehend long-context textual information, expanding the context window to accommodate even up to 200K input tokens. Meanwhile, building high-quality benchmarks with much longer text lengths and more demanding tasks to provide comprehensive evaluations is of immense practical interest to facilitate long context understanding research of LLMs. However, prior benchmarks create datasets that ostensibly cater to long-text comprehension by expanding the input of traditional tasks, which falls short to exhibit the unique characteristics of long-text understanding, including long dependency tasks and longer text length compatible with modern LLMs' context window size. In this paper, we introduce a benchmark for eXtremely Long context understanding with Long-range dependencies, **XL²Bench**, which includes three scenarios—Fiction Reading, Paper Reading, and Law Reading—and four tasks of increasing complexity: Memory Retrieval, Detailed Understanding, Overall Understanding, and Open-ended Generation, covering 27 subtasks in English and Chinese. It has an average length of 100K+ words (English) and 200K+ characters (Chinese). Evaluating six leading LLMs on XL²Bench, we find that their performance significantly lags behind human levels. Moreover, the observed decline in performance across both the original and enhanced datasets underscores the efficacy of our approach to mitigating data contamination.¹

1 Introduction

Large Language Models (LLMs) have attracted considerable interest for their remarkable capabilities

¹Code and benchmark are available at <https://github.com/anonymous/XL2Bench>

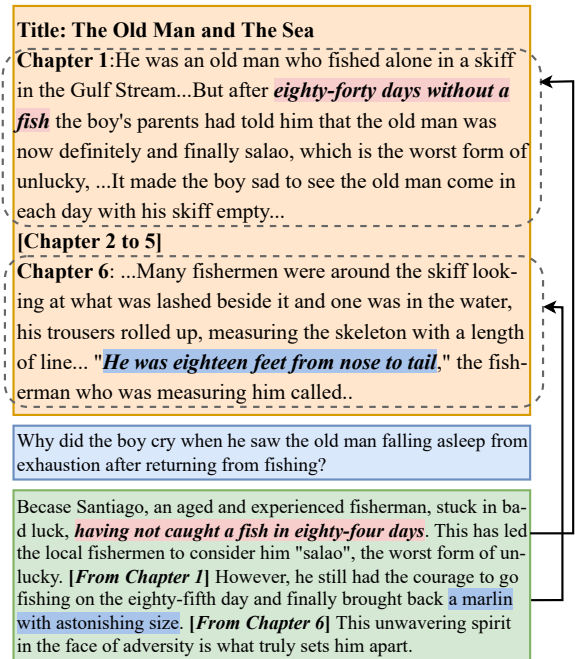


Figure 1: An illustrative example of long-dependency task, in which the model needs to make connective inferences across input document to fulfill the goal.

ities in a wide range of NLP tasks. However, a common limitation among these models is the fixed context window size (for example, LLaMA with maximum 2048 tokens and GPT-3.5 with maximum 4096 tokens), rendering them incapable of memorizing and understanding extremely long inputs (Liu et al., 2023). Evidenced by a basic passkey retrieval task, the accuracy of LLaMA recalling a passkey plummets from nearly 100% to nil when the text surpasses 2048 tokens (Tworkowski et al., 2023).

In pursuit of the goal of improving LLM's ability to comprehend long-context textual information, various efforts have been proposed to expand the context window of LLMs, such as sparse attention (Tworkowski et al., 2023; Chen et al., 2023; Mohtashami and Jaggi, 2023), length ex-

059 trapolation (Dai et al., 2019; Su et al., 2021; Peng
060 et al., 2023), and context compression (Ge et al.,
061 2023; Mu et al., 2023). Given the notable ad-
062 vances achieved by these techniques, the neces-
063 sity for high-quality benchmarks, featuring longer
064 text lengths and more complex tasks, is escalating
065 to facilitate thorough evaluations of LLMs’ long
066 context understanding ability.

067 Being able to understand long-range depen-
068 dencies in context and be sensitive to various
069 perturbations applied to distant context is what
070 sets long text understanding apart from traditional
071 NLP tasks (Wang et al., 2020; Tay et al., 2021;
072 Rae and Razavi, 2020). Existing benchmarks for
073 long-text understanding, such as LongBench (Bai
074 et al., 2023), L-Eval (An et al., 2023), and In-
075 finiteBench (Zhang et al., 2023c), often merely
076 expand the input of traditional tasks to create
077 datasets that ostensibly cater to long-text compre-
078 hension (Bai et al., 2023; An et al., 2023). However,
079 this approach does not tailor tasks to the distinct
080 features of long-text comprehension, thereby im-
081 peding the thorough assessment of LLMs’ abilities
082 in understanding extended contexts. Moreover, the
083 average text length in existing benchmarks usu-
084 ally does not exceed a few thousand tokens, sig-
085 nificantly shorter than the long texts perceived in
086 human cognition. For example, a user might up-
087 load an entire novel and inquire about the devel-
088 opment of the protagonist’s storyline. This task
089 would require the model to process and compre-
090 hend texts spanning over ten thousands of words,
091 necessitating long-range understanding and reason-
092 ing within the content to adequately address the
093 question. Traditional benchmarks typically fall
094 short in measuring capabilities of LLMs to ag-
095 gregate disparate pieces of information scattered
096 throughout the whole input texts in more realis-
097 tic scenarios, making it challenging to truly eval-
098 uate LLMs’ ability on long context understand-
099 ing (Dong et al., 2023; Kwan et al., 2023).

100 In light of the deficiencies identified in cur-
101 rent benchmarks, this paper proposes a bench-
102 mark for eXtremely Long context understanding
103 with Long-range dependencies, **XL²Bench**, which
104 features three scenarios—Fiction Reading, Pa-
105 per Reading, and Law Reading. XL²Bench con-
106 tains extremely long documents with an average
107 of 100K+ words (English) and 200K+ characters
108 (Chinese), along with 632K questions spanning
109 over four specifically designed tasks to examine

a model’s ability to aggregate and compare infor- 110
mation across long context, including *Memory Re-* 111
trieval, *Detailed Understanding*, *Overall Under-* 112
standing, and *Open-ended Generation*. These tasks 113
mimic the way people use LLMs in real-world sce- 114
narios. Figure 1 illustrates a case where the model 115
explains a boy’s tears as stemming from a story 116
about the old man who, against significant chal- 117
lenges, successfully captures a marlin. To construct 118
a solid answer, it demands the model to identifies 119
passages describing the boy’s reaction, the man’s 120
triumph, and his earlier hardships across various 121
chapters, and make connective inferences using 122
details buried far back in the long context. 123

124 Besides, to address data contamination caused 125
by outdated long texts contained in benchmark, we 126
implement three data augmentation strategies: **text** 127
transformation, which involves altering the origi- 128
nal text into a different language or style; **text re-** 129
placement, which entails modifying or substituting 130
key textual information; and **text concatenation**, 131
which incorporates integrating additional texts into 132
the original document. 133

134 Results of experiments on multiple state-of-the- 135
art LLMs reveal that even the most advanced LLMs 136
currently available fall short of reaching human- 137
level proficiency on XL²Bench. Despite these mod- 138
els’ ability to handle texts of considerable length, 139
there is a marked decline in performance as the 140
text lengthens. Additionally, the results obtained 141
by RAG (Li et al., 2022a; Gao et al., 2023) on 142
XL²Bench demonstrate that retrieval-based meth- 143
ods fail in overall and detailed understanding tasks; 144
instead, they require that the models comprehen- 145
sively grasp the entirety of the long texts. Further- 146
more, we conduct ablation experiments to com- 147
pare model performance on both original and aug- 148
mented benchmarks, which shows that the strate- 149
gies we employ to address the issue of data con- 150
tamination are indeed effective. 151

Our contributions are delineated as follows: 150

- We construct XL²Bench, a comprehensive 151
benchmark for extremely long text under- 152
standing with well-designed tasks. 153
- We formulate three data augmentation tech- 154
niques to circumvent the issue of data con- 155
tamination frequently encountered when us- 156
ing LLMs alongside existing NLP datasets. 157
Through experimentation, we validate the ef- 158
ficacy of these methodologies in mitigating 159
concerns about data contamination. 160

- We conduct empirical experiments to evaluate the performance of advanced LLMs using XL²Bench. The results reveal that contemporary LLMs are still facing challenges in achieving comprehensive understanding across long textual inputs.

2 Related Work

2.1 Long Context Modeling

Large language models (LLMs), such as GPT-4 (Achiam et al., 2023) and Llama (Touvron et al., 2023a,b), have exhibited superior performance across a variety of text generation tasks and practical deployment scenarios (Zhao et al., 2023; Wan et al., 2023; Guo et al., 2023). Nonetheless, the principal limitation hindering LLMs from harnessing their greater potential is the context window size—the upper limit of text length the model is capable of processing (Ratner et al., 2023). To circumvent this limitation, methods based on Position Encoding (Shaw et al., 2018), length extrapolation (Newman et al., 2020), and sparse attention mechanisms (Zhang et al., 2021; Gao and Liu, 2023), such as Alibi (Press et al., 2022), RoPE (Su et al., 2021), and Landmark (Mohtashami and Jaggi, 2023), have been presented. Furthermore, some strategies compress texts to align with the model’s context window size (Mu et al., 2023; Chevalier et al., 2023). Alternative approaches like Retrieval-Augmented Generation (Cai et al., 2022; Li et al., 2022b) and Memory Bank (Wang et al., 2023) utilize segmented retrieval followed by generation.

2.2 Evaluation Benchmarks

Current research is frequently directed at developing benchmarks tailored to specific tasks, such as reasoning (Li et al., 2023), code (Chen et al., 2021; Austin et al., 2021), and mathematics (Hendrycks et al., 2021; Cobbe et al., 2021; Zhang et al., 2023b). However, existing benchmarks, such as LongBench (Bai et al., 2023), L-Eval (An et al., 2023), and Bamboo (Dong et al., 2023), essentially expand existing NLU datasets, which may not pose sufficient difficulty and are prone to data contamination, and often fall short in text length. Besides, M4LE (Kwan et al., 2023) offers control over text length within benchmarks. It constructs texts from fragments of multiple summarization datasets, which compromises textual cohesion. InfiniteBench (Zhang et al., 2023c) introduces a

broader range of tasks. However, the manual annotation required for such a benchmark is extremely costly. By way of contrast, XL²Bench leverages LLMs and meticulous human review to construct the benchmark cost-effectively.

3 Methodology

In this section, we introduce the construction methodologies of XL²Bench and design of tasks with various level of difficulty.

3.1 Task Design

We evaluate the model’s understanding of extremely long texts from the perspectives of fine-grained retrieval and coarse-grained understanding. Based on this, we design four tasks: *Memory Retrieval*, *Detailed Understanding*, *Overall Understanding*, and *Open-ended Generation*.

Memory Retrieval. This task challenges the model to accurately retrieve and respond to queries by finding content within the text that aligns with given instructions. For instance, the model may be asked to pinpoint the specifics of a legal entry within a law or identify the originating chapter of a passage from a novel, thereby evaluating its capability to accurately locate and interpret question-relevant content.

Detailed Understanding. Here, the model is tasked with not only retrieving content but also comprehensively understanding it to perform activities such as summarization or question answering. This demands a more profound level of textual comprehension, surpassing mere content retrieval to include an in-depth analysis and synthesis of the text.

Overall Understanding. To circumvent tasks being completed through simple content retrieval, we introduce the Overall Understanding task. This task necessitates a holistic comprehension of the long text, enabling the model to build long-range dependencies and tackle inquiries related to overarching themes, such as the depiction of a character throughout a novel or the trajectory of a company’s stock across its history.

Open-ended Generation. Building on a solid foundation of long text understanding, the model is expected to undertake generation tasks rooted in it, such as role-playing a character in the fiction. Outputs should demonstrate creative expansion and

Tasks	Subtasks	Source	Num		Avg. Len		Metric
			CN	EN	CN	EN	
Fiction Reading							
Memory Retrieval	Content Location	Content Extraction	1495	1405	571.6K	111.5K	Acc.
	Content Retrieval	Content Extraction	299	261	571.1K	116.0K	Acc.
Detailed Understanding	Chapter Summarization	Data Synthesis	167	156	569.7K	110.6K	Rouge-L
	Question Answering	Data Synthesis	249	269	562.0K	114.7K	BLEU
Overall Understanding	Chapter Counting	Content Extraction	30	27	569.7K	113.4K	Acc.
	Background Summarization	Data Synthesis	30	27	570.3K	113.7K	Rouge-L
	Event Extraction	Data Synthesis	30	27	570.2K	113.7K	Rouge-L
	Fiction Summarization	Data Synthesis	30	27	570.4K	113.8K	Rouge-L
	Character Description	Data Synthesis	191	140	589.7K	143.5K	Rouge-L
Open-ended Generation	Relationship Analysis	Data Synthesis	193	432	606.3K	189.8K	Rouge-L
	Role-play Conversation	Data Synthesis	293	256	592.7K	115.2K	BLEU
	News Generation	Data Synthesis	30	27	570.7K	114.0K	BLEU
	Poem Generation	Data Synthesis	30	27	570.1K	113.6K	BLEU
Paper Reading							
Memory Retrieval	Content Retrieval	Content Extraction	-	4532	-	13.7K	Acc.
Detailed Understanding	Section Summarization	Data Synthesis	-	3136	-	14.1K	Rouge-L
	Terminology Explanation	Data Synthesis	-	14981	-	13.5K	BLEU
Overall Understanding	Paper Counting	Content Extraction	-	3100	-	13.5K	Acc.
	Paper Summarization	Data Integration	-	518	-	14.0K	Rouge-L
Open-ended Generation	Paper Review	Data Integration	-	518	-	14.0K	BLEU
	Rating Score	Data Integration	-	518	-	13.6K	MAE
Law Reading							
Memory Retrieval	Legal Entry Location	Content Extraction	2213	-	105.6K	-	Acc.
	Legal Entry Retrieval	Content Extraction	2225	-	105.3K	-	Acc.
Detailed Understanding	Legal Definition QA	Data Synthesis	2635	-	102.9K	-	BLEU
	Legal Number QA	Data Synthesis	1477	-	105.7K	-	Acc.
Overall Understanding	Legal Entry Counting	Content Extraction	122	-	103.0K	-	Acc.
	Multiple Choice QA	Data Integration	16881	-	95.6K	-	F1
Open-ended Generation	Case Adjudication	Data Integration	588369	-	72.7K	-	Acc.

Table 1: An overview of the statistics of XL²Bench. **Source** represents the method we use to construct the dataset for this subtask. **Num** represents the number of <input, output> pairs this subtask possesses. **Avg. Len** denotes the average combined length of the input and output, which is computed using the number of characters for Chinese and the number of words for English. **K** stands for 1024. For example, 200K = 200*1024.

inference, adhering to the text’s core themes and concepts, while ensuring originality and thematic consistency.

Table 1 delineates the various subtasks encapsulated within these four primary tasks. For more task descriptions of XL²Bench, please refer to Appendix A.

3.2 Benchmark Construction

In this subsection, we describe the sources from which we gather data and the methodologies we employ for constructing the benchmark in three different scenarios.

We gather long texts categorized under three scenarios. For fiction reading, we select a variety of novels written in both Chinese and English. For paper reading, we download PDF versions and reviews of papers submitted to ICLR 2023 from

Openreview². For law reading, we gather a substantial collection of original Chinese legislations.

To minimize cost of human annotation, we employ three methods to construct : *Content Extraction*, *Data Integration*, and *Data Synthesis*.

Content Extraction. We extract content from the original text to serve as the answer and use the index of this portion of the content to formulate the question. For instance, we used the title of a paper as the answer, with the corresponding question being: *What is the title of this paper?*

Data Integration. Tasks within certain short text datasets bear formal resemblance to what we have designed, exemplified by Document QA. Consequently, we contemplate leveraging these datasets

²<https://openreview.net/group?id=ICLR.cc/2023/Conference>

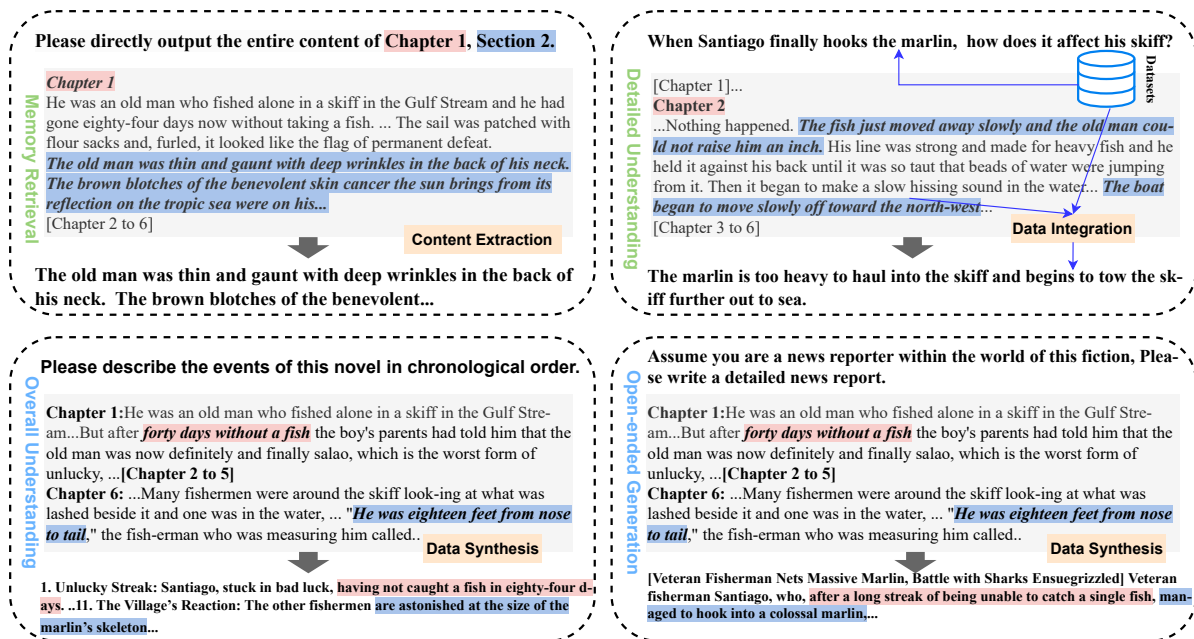


Figure 2: Illustration of the designed long context understanding tasks.

to augment our benchmark. More precisely, we employ LLMs to facilitate the alignment of data from the pre-existing datasets with our collected long texts. In an effort to mitigate potential overestimations of performance resulting from the model’s familiarity with these datasets during its training phase, we utilize LLMs to meticulously rewrite the original texts and remove any information that may indicate the data source.

Data Synthesis. For tasks that lack corresponding datasets, we utilize LLMs for direct generation. For summarization tasks, we employ structured text summarization method (Chang et al., 2023) via LLMs. For QA tasks, we use in-context learning (Brown et al., 2020) to construct some examples for the model to generate.

Employing the aforementioned approaches, we have constructed an extremely-long text benchmark encompassing three distinct scenarios, four overarching tasks, 27 detailed subtasks, and a corpus of 700+ texts with a average length of 100K+ words for English and 200K+ characters for Chinese. The statistics of our benchmark are shown in Table 1.

3.3 Data Contamination

The potential of data contamination warrants serious consideration when constructing a benchmark (Sainz et al., 2023; Deng et al., 2023; Magar and Schwartz, 2022). The risk arises when the test set data is either identical to, or strikingly similar to, the training set data. This could result in the model

memorizing specific answers instead of acquiring the ability to reason or generalize from unseen data. In our construction process, the selected novels, academic papers, and legal texts may have been included in the training corpus of LLMs. Consequently, the model may not need to fully comprehend the entire text to accomplish various tasks. In order to mitigate the impact of data contamination on model’s performance, we follow Yang et al. (2023) and adopt three strategies, namely *text transformation*, *key information replacement*, and *text concatenation* for fiction data augmentation.

Text Transformation. We utilize LLMs to facilitate mutual translation of fictions between Chinese and English, whereby the original Chinese (English) novels are rendered into English (Chinese). In accordance, the input and output for each task are also translated.

Key Information Replacement. We employ LLMs to extract key information from a chapter or section, such as names, places, and times. We then generate corresponding texts to replace these elements, resulting in a collection of ⟨original text - replacement text⟩ pairs, which are subsequently used for content substitution throughout the entire text and tasks.

Text Concatenation. We insert a short story into the original fiction as one of its chapters, and use this template to bridge: *Now, let’s pause the current story narration and turn to a new story*[New

Models	MR		DU		OU						TG		
	C-L	C-R	C-S	QA	C-C	B-S	E-E	F-S	Ch-D	Re-A	RP-C	N-G	P-G
YaRN-Mistral-7B	<1	<1	4.46	2.26	13.78	8.09	16.17	5.52	8.35	7.91	7.28	4.42	5.91
InternLM2-C-7B	<1	<1	8.27	<1	6.67	11.68	9.97	11.97	6.92	2.22	1.16	5.88	3.49
InternLM2-C-20B	6.85	<1	17.22	9.82	53.33	15.58	18.61	17.29	21.98	28.92	11.65	16.67	10.09
Kimi-Chat	<u>60.39</u>	<u>17.23</u>	23.53	<u>33.13</u>	86.30	24.32	20.08	25.10	<u>22.24</u>	54.99	12.81	<u>27.31</u>	<u>12.22</u>
GLM-4	63.44	20.08	18.12	14.51	<u>72.73</u>	18.40	<u>20.42</u>	15.84	<u>22.22</u>	42.27	<u>13.62</u>	19.70	11.69
GPT-4-Turbo	54.36	11.89	<u>19.87</u>	37.23	<u>60.00</u>	<u>21.21</u>	21.40	<u>21.57</u>	23.14	<u>49.05</u>	17.58	30.19	16.56

Table 2: Results (%) of six LLMs on Chinese Fiction Reading. **MR, DU, OU, TG** are the abbreviations for the initials of four tasks. **C-L, C-R, C-S**, etc., represent the abbreviations of 13 subtasks. The context window size of GPT-4-Turbo and YaRN-Mistral-7B is 128K, whereas it is 200K for other models. The **bold** numbers in the results represent the best scores, whereas the underlined numbers indicate the second-best scores.

Story]The story is over, let’s get back to the original fiction. Then, we merge the data in four tasks of this short story with the original fiction.

Through above three strategies, we construct **Fiction-T** (Translated), **Fiction-R** (Replaced), and **Fiction-C** (Concatenated). These three datasets can ensure that the model must fully comprehend the entire text in order to accomplish tasks, rather than being able to complete tasks by recalling the content of the training phase.

3.4 Implementation Details

We select GPT-4-Turbo (Achiam et al., 2023) to help us construct XL²Bench. GPT-4 currently stands as the highest-performing LLMs, characterized by a 128k context window along with superior memory, reasoning, and generation capabilities. To ensure optimum quality of the benchmark, we enlist the assistance of several university students to manually review the content generated by GPT-4-Turbo. The prompts and input templates used throughout the construction process are available in our GitHub repository due to space limit.

4 Experimental Settings

4.1 Generative Large Language Models

We introduce current LLMs with context window size **more than 100k** evaluated in our experiments. Models such as LLama2 (Touvron et al., 2023b) and ChatGLM2 (Zeng et al., 2023) have context window size significantly shorter than the average text length of XL²Bench, resulting in an excessive need to truncate texts, which leads to suboptimal performance. Consequently, we do not evaluate the effectiveness of these models.

GPT-4-Turbo Developed by OpenAI, GPT-4-Turbo represents the pinnacle of current advancements, demonstrating exceptional reasoning and

instruction-following capacities. It is distinguished by its extensive context window of 128K tokens. We employ this model via API³.

GLM-4 GLM-4 is the latest model developed by Zhipu AI. Compared to ChatGLM2, it boasts more powerful question-answering and text generation capabilities, capable of processing up to 200,000 tokens. We employ this model via API⁴.

Kimi Chat Kimi Chat, developed by Moonshot AI, boasts exceptional performance in processing extremely-long text inputs of up to 200K tokens. We employ this model via API⁵.

InternLM2-Chat Equipped with a 200k context window, InternLM2 exhibits comprehensive enhancements across all functionalities when juxtaposed with the previous generation model. We employ InternLM2-Chat-7B-200k and InternLM2-Chat-20B-200k.

YaRN-Mistral The computationally efficient length extrapolation technology YaRN makes it possible to expand LLM’s context window size while conserving resources. We leverage YaRN-Mistral-7B-128k.

4.2 Retrieval-Augmented Generation Methods

One type of methods to handle long texts with small context window size in LLMs is Retrieval-Augmented Generation (RAG) (Li et al., 2022a). Given a long context, we first splits it into chunks. Then, using a specific retriever, we compute the embedding of the text chunks and query. Only the top-N chunks, based on the cosine similarity of their embeddings to the query embedding,

³<https://chat.openai.com/>

⁴<https://open.bigmodel.cn/>

⁵<https://www.moonshot.cn/>

Models	MR		DU		OU		TG
	LE-L	LE-R	Def-QA	Num-QA	LE-C	MCQA	Case-Adj
YaRN-Mistral-7B-128K	<1	11.29	8.62	<1	3.36	<1	<1
InternLM2-Chat-7B-200K	<1	2.61	3.52	<1	<1	<1	<1
InternLM2-Chat-20B-200K	5.41	22.60	40.57	58.03	11.76	44.23	41.05
Kimi-Chat-200K	32.61	88.83	48.08	<u>63.85</u>	28.10	<u>63.11</u>	<u>47.40</u>
GLM-4-200K	<u>16.97</u>	<u>72.76</u>	<u>43.17</u>	67.63	31.14	53.56	47.31
GPT-4-Turbo-128K	13.41	63.48	40.26	62.50	<u>29.51</u>	63.24	48.89

Table 3: Results (%) of six LLMs on Law Reading. **LE-L**, **LE-R**, **Def-QA**, **Num-QA**, **LE-C**, **MCQA** and **Case-Adj** represent *Legal Entry Location*, *Legal Entry Retrieval*, *Legal Definition QA*, *Legal Number QA*, *Legal Entry Counting*, *Multiple Choice QA* and *Case Adjudication*, respectively. Rest settings remain the same as in the previous tables.

Models	MR		DU		OU		TG
	LE-L	LE-R	Def-QA	Num-QA	LE-C	MCQA	Case-Adj
InternLM2-Chat-20B-200K	5.41	22.60	40.57	58.03	11.76	44.23	41.05
w/ Sentence-Transformers	<1	16.54	11.59	11.22	4.92	39.92	31.16
w/ LLM-Embedder	1.86	21.68	11.97	19.98	2.46	42.59	38.83
w/ Contriever	<1	16.73	10.23	5.44	4.10	40.23	37.79

Table 4: Results (%) of InternLM2-Chat-20B-200K using different embedding models on Law Reading. *w/* represents *with*. The best performance over of each subtask is in **bold**.

are concatenated. These top-N chunks along with the query are then fed into the model to produce an answer. We test this technique’s impact on LLMs evaluation results, to see if the model could complete XL²Bench tasks by retrieving certain fixed chunks. We employ LangChain⁶ and three retrievers: Sentence-Transformers (Reimers and Gurevych, 2020), LLM-Embedder (Zhang et al., 2023a), and Contriever (Izacard et al., 2022). We set the chunk size to 500 and N=5.

4.3 Automatic Evaluation Metrics

For tasks with fixed answers, such as Content Location in Fiction Reading, we adopt **Accuracy** as an intuitive measure to demonstrate the model’s performance. For MCQA, we utilize **F1-Score** to objectively evaluate the model’s capability to accurately answer all the correct options. For summary tasks, we select **Rouge-L** to reflect whether the model can correctly identify key information in a document. For generative tasks, we employ **BLEU** to measure the congruence between the generated content by model and the reference content. For Rating Score subtask, we choose MAE to calculate the average absolute difference between predicted and true scores. Details can be found in Table 1.

⁶https://python.langchain.com/docs/get_started/introduction

4.4 Inference Settings

We conduct the evaluation in a zero-shot setting. The input templates we use during inference can be found in Appendix B. When the input length exceeds the context window size of LLMs, we truncate the input sequence from the middle, as the front and end of the sequence may contain crucial information such as instructions or questions. For models that are API-callable, we follow the original settings provided in the sample code of these models. For locally deployed models, we select the decoding parameters as follows: Temperature=0.2, Top-K=40, Top-P=0.9, Repetition Penalty=1.02.

5 Results and Analysis

5.1 Long Texts Processing

The results pertaining to three scenarios are delineated in Table 2 and 3. Due to space constraints, the remaining results are relegated to Appendix C. The key findings from the experiments can be summarized below.

The overall performance of all LLMs is notably unsatisfactory. Regardless of whether they are open-source or closed-source, LLMs consistently score low across various metrics pertaining to the 27 subtasks, particularly in retrieval and counting tasks where human performance approaches

100%. We hypothesize that these results are attributable to the use of sparse attention or length extrapolation techniques within the extended model context window, as well as the truncation operation employed when the input text is too long.

Closed-source models outperform open-source models. The comparative performance analysis of three closed-source LLMs demonstrates a superior performance over their open-source counterparts. Furthermore, with 7B parameters, YaRN-Mistral and InternLM2-Chat-7B exhibit sub-optimal performance across a majority of tasks, achieving scores below 1. This demonstrates the importance of the model’s parameter size for effectively managing tasks in XL²Bench.

LLMs have a preference for the language of the input text. GLM-4 and Kimi-Chat performs well on Chinese-language tasks (Law Reading and Fiction-CN), while GPT-4 performs well on English-language tasks (Paper Reading and Fiction-EN). We infer that this may be due to the different proportions of Chinese and English datasets used in the training process of these three models. This further indicates that the dataset is a particularly critical factor that affects model performance.

GPT-4’s performance on self-generated sub-tasks does not meet expectations. In particular, for subtasks where the ground truth is established by GPT-4 itself, we meticulously assessed the model’s efficacy. Contrary to our initial assumptions, GPT-4’s scores on these tasks were lower than anticipated. Upon an in-depth analysis of the model-generated content, we hypothesized that the verbose nature of the text could have adversely affected GPT-4’s understanding of the task descriptions, leading to a diminished output quality.

The findings and analyses presented above indicate that existing context window expansion technologies fall significantly short of reaching or approximating human-level performance. Addressing the issue of context dependency represents a critical area for potential breakthroughs and merits further exploration.

5.2 Performance of Retrieval-Augmented Generation Methods

In this subsection, we assess the performance of InternLM2-Chat-20B-200K, which utilizes three distinct retrievers on Law Reading scenarios. Results illustrated in Table 4, indicate a uniform reduction in the model’s performance across all subtasks

following the adoption of RAG methods. Notably, the most substantial declines in performance were observed in the Definition QA and Number QA tasks. We postulate that these decreases may be due to the retrievers’ failure to recall relevant segments of text. The results and subsequent analysis imply that effectively addressing the tasks in XL²Bench demands more than merely retrieving relevant documents.

5.3 Impact of Context Length

In this subsection, we explore the impact of context length on the performance of LLMs. Our evaluation focuses on the average performance of the InternLM2-Chat-20B across four tasks, using legal texts of varying lengths. Results presented in Appendix D illustrate that the model’s performance significantly declines with longer texts, as evidenced by a steeper curve. This observation underscores the model’s challenges in effectively managing the complexities of long text modeling.

5.4 Impact of Data Contamination

In this subsection, we conduct an ablation study to examine the effectiveness of the methodologies employed to reduce data contamination. The results indicate that our data augmentation techniques can, to some extent, reduce the likelihood of biased evaluations. A detailed discussion is provided in Appendix E due to space limit.

6 Conclusion

In this paper, we present XL²Bench, a comprehensive benchmark for extremely long text understanding with long-range dependencies. XL²Bench consists of three scenarios, four tasks, and 27 sub-tasks, with an average length of over 100K words (English) and 200K characters (Chinese). We automatically construct the benchmark via LLMs, significantly reducing the cost of manually annotating the datasets. Furthermore, we mitigate data contamination risks through carefully designed techniques. Extensive experiments on XL²Bench yield insights into the capabilities of current LLMs for long text understanding. We also demonstrate that RAG methods are not suitable for XL²Bench as the benchmark requires a comprehensive understanding of the entire text to complete the tasks. Results and analyses indicate that XL²Bench is a valuable resource for advancing research in the comprehension of long texts.

567
568
569
570
571
572
573
574
575
576

577

578
579
580
581
582
583
584
585

586

587
588
589
590
591

592
593
594
595

596
597
598
599
600

601
602
603
604
605
606

607
608
609
610
611
612

613
614
615
616
617

Limitations

The limitations of XL²Bench mainly come from the disadvantages of using LLMs. First of all, most of the large language models that work well are not open source or free. This makes it difficult to conduct batch experiments or daily use on it. Next, a small number of open-source models require a lot of GPU resources when used, which is a difficult problem for quite many researchers, such as students.

Ethics Statement

We honor and support the ACL code of Ethics. Our benchmark XL²Bench aims to evaluate large language models' ability of long-text comprehension. The interaction and assistance process do not involve any bias towards to the participants. Following our thorough examination, we can confirm that our benchmark is free from any privacy or ethical concerns.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Chenxin An, Shansan Gong, Ming Zhong, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. 2023. L-eval: Instituting standardized evaluation for long context language models. *CoRR*, abs/2307.11088.

Jacob Austin, Augustus Odena, Maxwell I. Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. 2021. Program synthesis with large language models. *CoRR*, abs/2108.07732.

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2023. Longbench: A bilingual, multitask benchmark for long context understanding. *CoRR*, abs/2308.14508.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Deng Cai, Yan Wang, Lemao Liu, and Shuming Shi. 2022. Recent advances in retrieval-augmented text generation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3417–3419.

Yapei Chang, Kyle Lo, Tanya Goyal, and Mohit Iyyer. 2023. Boookscore: A systematic exploration of book-length summarization in the era of llms. *arXiv preprint arXiv:2310.00785*. 618
619
620
621

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code. *CoRR*, abs/2107.03374. 622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642

Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. 2023. Longlora: Efficient fine-tuning of long-context large language models. *CoRR*, abs/2309.12307. 643
644
645
646

Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. 2023. Adapting language models to compress contexts. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 3829–3846. Association for Computational Linguistics. 647
648
649
650
651
652
653

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168. 654
655
656
657
658
659

Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc Viet Le, and Ruslan Salakhutdinov. 2019. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 2978–2988. Association for Computational Linguistics. 660
661
662
663
664
665
666
667
668

Chunyan Deng, Yilun Zhao, Xiangru Tang, Mark Gestein, and Arman Cohan. 2023. Investigating data contamination in modern benchmarks for large language models. *CoRR*, abs/2311.09783. 669
670
671
672

Zican Dong, Tianyi Tang, Junyi Li, Wayne Xin Zhao, and Ji-Rong Wen. 2023. BAMBOO: A comprehensive benchmark for evaluating long text mod- 673
674
675

676	eling capacities of large language models. <i>CoRR</i> , abs/2309.13345.	728
677		729
678	Yue Gao and Jian-Wei Liu. 2023. Adaptively sparse transformers hawkes process . <i>Int. J. Uncertain. Fuzziness Knowl. Based Syst.</i> , 31(4):669–689.	730
679		731
680		732
681	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey . <i>CoRR</i> , abs/2312.10997.	733
682		734
683		735
684		736
685		737
686	Tao Ge, Jing Hu, Xun Wang, Si-Qing Chen, and Furu Wei. 2023. In-context autoencoder for context compression in a large language model . <i>CoRR</i> , abs/2307.06945.	738
687		739
688		740
689		741
690	Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, and Deyi Xiong. 2023. Evaluating large language models: A comprehensive survey . <i>CoRR</i> , abs/2310.19736.	742
691		743
692		744
693		745
694		746
695	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	747
696		748
697		749
698		750
699		751
700		752
701		753
702	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning . <i>Trans. Mach. Learn. Res.</i> , 2022.	754
703		755
704		756
705		757
706		758
707	Wai-Chung Kwan, Xingshan Zeng, Yufei Wang, Yusen Sun, Liangyou Li, Lifeng Shang, Qun Liu, and Kam-Fai Wong. 2023. M4LE: A multi-ability multi-range multi-task multi-domain long-context evaluation benchmark for large language models . <i>CoRR</i> , abs/2310.19240.	759
708		760
709		761
710		762
711		763
712		764
713	Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023. CMMLU: measuring massive mult-task language understanding in chinese . <i>CoRR</i> , abs/2306.09212.	765
714		766
715		767
716		768
717		769
718	Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. 2022a. A survey on retrieval-augmented text generation . <i>CoRR</i> , abs/2202.01110.	770
719		771
720		772
721	Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. 2022b. A survey on retrieval-augmented text generation . <i>arXiv preprint arXiv:2202.01110</i> .	773
722		774
723		775
724	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts . <i>CoRR</i> , abs/2307.03172.	776
725		777
726		778
727		779
	Inbal Magar and Roy Schwartz. 2022. Data contamination: From memorization to exploitation . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022</i> , pages 157–165. Association for Computational Linguistics.	780
		781
		782
	Amirkeivan Mohtashami and Martin Jaggi. 2023. Landmark attention: Random-access infinite context length for transformers . <i>CoRR</i> , abs/2305.16300.	783
		784
	Jesse Mu, Xiang Lisa Li, and Noah D. Goodman. 2023. Learning to compress prompts with gist tokens . <i>CoRR</i> , abs/2304.08467.	785
		786
	Benjamin Newman, John Hewitt, Percy Liang, and Christopher D. Manning. 2020. The EOS decision and length extrapolation . In <i>Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP, BlackboxNLP@EMNLP 2020, Online, November 2020</i> , pages 276–291. Association for Computational Linguistics.	787
		788
	Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023. Yarn: Efficient context window extension of large language models . <i>CoRR</i> , abs/2309.00071.	789
		790
	Ofir Press, Noah A. Smith, and Mike Lewis. 2022. Train short, test long: Attention with linear biases enables input length extrapolation . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	791
		792
	Jack W. Rae and Ali Razavi. 2020. Do transformers need deep long-range memory? In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 7524–7529. Association for Computational Linguistics.	793
		794
		795
	Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram, Inbal Magar, Omri Abend, Ehud Karpas, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. Parallel context windows for large language models . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 6383–6402. Association for Computational Linguistics.	796
		797
	Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	798
		799
	Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 10776–10787. Association for Computational Linguistics.	800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827

786	Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018.	Shuohang Wang, Luowei Zhou, Zhe Gan, Yen-Chun	843
787	Self-attention with relative position representations.	Chen, Yuwei Fang, Siqi Sun, Yu Cheng, and Jingjing	844
788	In <i>Proceedings of the 2018 Conference of the North</i>	Liu. 2020. Cluster-former: Clustering-based sparse	845
789	<i>American Chapter of the Association for Computa-</i>	transformer for long-range dependency encoding.	846
790	<i>tional Linguistics: Human Language Technologies,</i>	<i>CoRR</i> , abs/2009.06097.	847
791	<i>NAACL-HLT, New Orleans, Louisiana, USA, June</i>	Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu,	848
792	<i>1-6, 2018, Volume 2 (Short Papers)</i> , pages 464–468.	Xifeng Yan, Jianfeng Gao, and Furu Wei. 2023. Aug-	849
793	Association for Computational Linguistics.	menting language models with long-term memory.	850
794	Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng	<i>CoRR</i> , abs/2306.07174.	851
795	Liu. 2021. Roformer: Enhanced transformer with	Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E.	852
796	rotary position embedding. <i>CoRR</i> , abs/2104.09864.	Gonzalez, and Ion Stoica. 2023. Rethinking bench-	853
797	Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen,	mark and contamination for language models with	854
798	Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang,	rephrased samples. <i>CoRR</i> , abs/2311.04850.	855
799	Sebastian Ruder, and Donald Metzler. 2021. Long	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang,	856
800	range arena : A benchmark for efficient transformers.	Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,	857
801	In <i>9th International Conference on Learning Repre-</i>	Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma,	858
802	<i>sentations, ICLR 2021, Virtual Event, Austria, May</i>	Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan	859
803	<i>3-7, 2021.</i> OpenReview.net.	Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023.	860
804	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	GLM-130B: an open bilingual pre-trained model. In	861
805	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	<i>The Eleventh International Conference on Learning</i>	862
806	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	<i>Representations, ICLR 2023, Kigali, Rwanda, May</i>	863
807	Azhar, Aurélien Rodriguez, Armand Joulin, Edouard	<i>1-5, 2023.</i> OpenReview.net.	864
808	Grave, and Guillaume Lample. 2023a. Llama: Open	Biao Zhang, Ivan Titov, and Rico Sennrich. 2021.	865
809	and efficient foundation language models. <i>CoRR</i> ,	Sparse attention with linear units. In <i>Proceedings</i>	866
810	abs/2302.13971.	<i>of the 2021 Conference on Empirical Methods in</i>	867
811	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	<i>Natural Language Processing, EMNLP 2021, Vir-</i>	868
812	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	<i>tual Event / Punta Cana, Dominican Republic, 7-11</i>	869
813	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	<i>November, 2021</i> , pages 6507–6520. Association for	870
814	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-	Computational Linguistics.	871
815	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	Peitian Zhang, Shitao Xiao, Zheng Liu, Zhicheng Dou,	872
816	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	and Jian-Yun Nie. 2023a. Retrieve anything to aug-	873
817	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	ment large language models.	874
818	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	Xiaotian Zhang, Chunyang Li, Yi Zong, Zhengyu Ying,	875
819	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	Liang He, and Xipeng Qiu. 2023b. Evaluating the	876
820	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	performance of large language models on GAOKAO	877
821	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	benchmark. <i>CoRR</i> , abs/2305.12474.	878
822	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	Xinrong Zhang, Yingfa Chen, Shengding Hu, Qihao Wu,	879
823	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	Junhao Chen, Zihang Xu, Zhenning Dai, Xu Han,	880
824	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2023c.	881
825	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	Infinitebench: 128k long-context benchmark for lan-	882
826	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	guage models.	883
827	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	884
828	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	Xiaolei Wang, Yupeng Hou, Yingqian Min, Be-	885
829	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	ichen Zhang, Junjie Zhang, Zican Dong, Yifan Du,	886
830	Melanie Kambadur, Sharan Narang, Aurélien Ro-	Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao	887
831	driguez, Robert Stojnic, Sergey Edunov, and Thomas	Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang	888
832	Scialom. 2023b. Llama 2: Open foundation and	Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen.	889
833	fine-tuned chat models. <i>CoRR</i> , abs/2307.09288.	2023. A survey of large language models. <i>CoRR</i> ,	890
834	Szymon Tworkowski, Konrad Staniszewski, Mikolaj	abs/2303.18223.	891
835	Pacek, Yuhuai Wu, Henryk Michalewski, and Piotr	A Task Descriptions	892
836	Milos. 2023. Focused transformer: Contrastive train-	In this section, we provide detailed descriptions of	893
837	ing for context scaling. <i>CoRR</i> , abs/2307.03170.	the input and output content of 27 subtasks. Please	894
838	Zhongwei Wan, Xin Wang, Che Liu, Samiul Alam,	note that the input includes a long text and an in-	895
839	Yu Zheng, Jiachen Liu, Zhongnan Qu, Shen Yan,	struction. We only describe the instruction.	896
840	Yi Zhu, Quanlu Zhang, Mosharaf Chowdhury, and		
841	Mi Zhang. 2023. Efficient large language models: A		
842	survey. <i>CoRR</i> , abs/2312.03863.		

897	A.1 Fiction Reading		
898	Content Location	Given the content of the fiction, the model outputs the location.	940
899			941
900	Content Retrieval	Given a location, the model outputs the corresponding fiction content.	942
901			
902	Chapter Summarization	Given a chapter number of the fiction, the model summarizes the corresponding chapter.	943
903			944
904			945
905	Question Answering	Give a detailed question about the fiction, the model outputs the answer.	946
906			947
907	Chapter Counting	The model outputs the quantity of the fiction.	948
908			949
909	Background Summarization	The model outputs the time background, place background, and social and cultural background of the fiction.	950
910			951
911			952
912	Event Extraction	The model outputs the main events of the fiction in chronological order.	953
913			954
914	Fiction Summarization	The model summarizes the whole fiction.	955
915			956
916	Character Description	The model outputs the description of the character in the fiction, including personality traits and personal experiences.	957
917			958
918			959
919	Relationship Analysis	The model outputs the relationship between two characters.	960
920			961
921	Role-play Conversation	Given a question, the model needs to assume the role of a character from the fiction to provide an answer.	962
922			963
923			964
924	News Generation	The model assume a news reporter within the world of the fiction, and reports on the final event involving the protagonist’s team, including the background of the event, the actions of the protagonist, the outcome, and the impact of the event.	965
925			966
926			967
927			968
928			969
929			
930	Poem Generation	The model writes a poem based on the core theme, key plot, important characters and specific context of the fiction.	970
931			971
932			972
933	A.2 Paper Reading		973
934	Content Retrieval	Given a location, the model outputs the corresponding paper content, such as title, authors.	974
935			975
936			976
937	Section Summarization	Given a section number of the paper, the model summarizes the corresponding section.	977
938			978
939			979
			980
			981
	Terminology Explanation	Given an scientific noun in the paper, the model outputs its explanation.	940
			941
			942
	Paper Counting	The model output the quantity of titles, authors, references, tables, figures, etc. of the paper.	943
			944
			945
	Paper Summarization	The model summarizes the whole paper.	946
			947
	Paper Review	The model assumes the role of a peer reviewer for an academic journal, and outputs a review of the paper, including: strengths and weaknesses.	948
			949
			950
			951
	Rating Score	The model assumes the role of a peer reviewer for an academic journal, and outputs a rating score of the paper from 0 to 10.	952
			953
			954
	A.3 Law Reading		955
	Legal Entry Location	Given the content of the law, the model outputs its corresponding index.	956
			957
	Legal Entry Retrieval	Given a locating of a legal entry, the mode outputs its content.	958
			959
	Legal Definition QA	Given a question about the law’s definitions, the model outputs the answer.	960
			961
	Legal Number QA	Given questions about the numbers in law, the model outputs the answer.	962
			963
	Legal Entry Counting	The model outputs the quantity of legal entries in this law.	964
			965
	Multiple Choices QA	Given a question with multiple choices, the model outputs the answer.	966
			967
	Case Adjudication	Given a legal case, the model outputs the verdict.	968
			969
	B Evaluation Input Templates		970
		For all texts and corresponding questions in XL ² Bench, we use the following template: <i>Please read the following text, and answer related question: [text] Question: [question] Directly output your answer without any additional analysis or explanation.</i>	971
			972
			973
			974
			975
			976
	C Results on English Fiction Reading and Paper Reading		977
			978
		We show the remaining results of six LLMs on English Fiction Reading and Paper Reading in Table 5 and Table 6.	979
			980
			981

Models	MR		DU		OU						TG		
	C-L	C-R	C-S	QA	C-C	B-S	E-E	F-S	Ch-D	Re-A	RP-C	N-G	P-G
YaRN-Mistral-7B	<1	<1	6.64	2.29	5.52	10.16	2.85	3.13	10.09	8.52	4.36	4.42	5.40
InternLM2-C-7B	<1	<1	3.08	<1	<1	7.73	5.15	4.57	7.01	2.31	6.90	4.23	21.88
InternLM2-C-20B	18.85	1.58	17.60	35.43	56.01	17.47	29.81	25.04	19.97	20.73	53.14	29.79	44.81
Kimi-Chat	38.19	33.56	24.46	<u>34.14</u>	88.89	30.30	38.79	39.16	<u>28.45</u>	25.46	37.10	<u>61.76</u>	<u>62.47</u>
GLM-4	26.68	<u>34.60</u>	18.06	<u>32.86</u>	66.67	28.75	34.46	24.30	25.24	27.56	<u>39.20</u>	<u>35.07</u>	<u>53.12</u>
GPT-4-Turbo	55.46	42.70	19.76	50.81	<u>77.50</u>	<u>29.30</u>	44.20	42.57	30.87	<u>27.16</u>	66.71	74.59	67.80

Table 5: Results (%) of six LLMs on English Fiction Reading.

Models	MR	DU		OU		TG	
	C-R	Sec-Sum	T-E	Paper-C	Paper-Sum	P-Review	R-Score↓
YaRN-Mistral-7B-128K	<1	10.19	15.86	11.69	5.04	33.23	None
InternLM2-Chat-7B-200K	<1	6.82	5.04	<1	7.31	39.80	None
InternLM2-Chat-20B-200K	25.84	24.91	30.27	33.37	34.41	45.11	<u>2.30</u>
Kimi-Chat-128k	<u>31.02</u>	<u>45.78</u>	31.43	44.44	36.68	66.04	4.39
GLM-4-200K	25.76	29.66	<u>33.40</u>	<u>47.62</u>	<u>36.91</u>	55.62	2.23
GPT-4-Turbo-200K	45.28	51.57	55.91	55.56	45.91	<u>62.12</u>	2.63

Table 6: Results of six LLMs on Paper Reading. **Sec-Sum**, **T-E**, **Paper-C**, **Paper-Sum**, **P-Review**, and **R-Score** represent *Section Summarization*, *Terminology Explanation*, *Paper Counting*, *Paper Summarization*, *Paper Review*, and *Rating Score* respectively. **None** signifies the model’s inability to generate a rating score, thus rendering it incapable of fulfilling the requirements of this subtask.

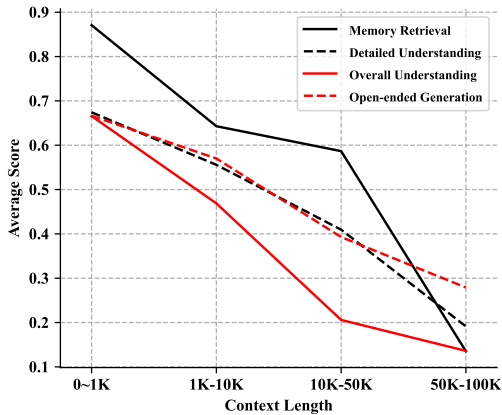


Figure 3: Average score (%) of four tasks under different context length on Law Reading.

D Impact of Context Length

Figure 3 illustrates that the model’s performance significantly declines with longer texts, as evidenced by a steeper curve. This observation underscores the model’s challenges in effectively managing the complexities of long text modeling.

E Ablation Study

In this section, we assess the effectiveness of our data augmentation strategies in mitigating the im-

991
992
993
994
995
996
997
998
999
1000

part of data contamination on model evaluation outcomes. We specifically examine the performance of the InternLM2-Chat-20B across different subsets of fiction data, namely Fiction, Fiction-T, Fiction-R, and Fiction-C, with the results detailed in Table 7. The observed reduction in performance across almost all subtasks within the augmented dataset indicates that our data augmentation techniques can, to some extent, reduce the likelihood of biased evaluations.

Scenarios	MR		DU		OU						TG		
	<i>C-L</i>	<i>C-R</i>	<i>C-S</i>	<i>QA</i>	<i>C-Q</i>	<i>F-B</i>	<i>F-E</i>	<i>F-S</i>	<i>Ch-D</i>	<i>Ch-R</i>	<i>Ch-DG</i>	<i>N-G</i>	<i>P-G</i>
Fiction	6.85	<1	17.22	9.82	53.33	15.58	18.61	17.29	21.98	28.92	11.65	16.67	10.09
Fiction-T	6.54	<1	12.28	5.05	52.16	10.21	10.80	6.67	2.28	13.89	12.36	11.89	5.01
Fiction-R	6.76	<1	5.11	6.48	53.33	8.04	11.72	4.96	3.33	17.67	12.12	11.84	5.78
Fiction-C	6.28	<1	5.23	3.39	53.33	7.65	4.46	13.41	2.49	15.56	13.79	12.68	7.91

Table 7: Results of six LLMs on Fiction, Fiction-T, Fiction-R, and Fiction-C.