

# Marathon: A Race Through the Realm of Long Context with Large Language Models

Lei Zhang<sup>1,2</sup> Yunshui Li<sup>1,2</sup> Ziqiang Liu<sup>1,2</sup> Jiayi Yang<sup>1,2</sup> Junhao Liu<sup>3</sup> Min Yang<sup>1</sup>

<sup>1</sup>Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences

<sup>2</sup>University of Chinese Academy of Sciences

<sup>3</sup>University of California, Irvine

{lei.zhang2, min.yang}@siat.ac.cn

## Abstract

Although there are currently many benchmarks available for evaluating the long context understanding and reasoning capability of large language models, with the expansion of the context window in these models, the existing long context benchmarks are no longer sufficient for evaluating the long context understanding and reasoning capability of large language models. In this paper, we have developed a fresh long context evaluation benchmark, which we name it **Marathon** in the form of multiple choice questions, inspired by benchmarks such as MMLU, for assessing the long context comprehension capability of large language models quickly, accurately, and objectively. We have evaluated several of the latest and most popular large language models, as well as three recent and effective long context optimization methods, on our benchmark. This showcases the long context reasoning and comprehension capabilities of these large language models and validates the effectiveness of these optimization methods. **Marathon** is available at <https://huggingface.co/datasets/LemonCoke/Marathon>.<sup>1</sup>

## 1 Introduction

With the rapid development of artificial intelligence technology, large language models represented by GPT-4 (OpenAI, 2023b) have demonstrated powerful capabilities. In addition to the well-known ChatGPT, the application of large language models has gradually extended to various industries. The large model in the medical field, MedGPT, can help people consult medical-related questions, while the large model in the legal field, LawGPT (Nguyen, 2023), can assist people in legal consultations, significantly reducing the barriers and costs for people to obtain knowledge. The large model in the programming field, Github Copilot (Chen

et al., 2021), can help developers complete their work faster, greatly improving productivity.

However, current large language models can only provide answers to users' questions and requests within a relatively small context window. There are two main reasons for this limitation. Firstly, large language models still lack the ability to extract crucial information from long contexts. As described in Liu et al. (2023) large language models can only effectively retrieve key information when it is located at the beginning or end of the prompt. When the key information is in the middle of the prompt, it becomes difficult to retrieve successfully. This phenomenon becomes more evident as the context length increases. Secondly, longer contexts require more computational resources and time for large language models to provide feedback, which not only increases the hardware requirements but also reduces the user experience. Therefore, the ability of large language models to handle long contexts becomes increasingly important, and comprehensive evaluation of their performance in this regard is crucial. Existing long context benchmarks, such as Longbench (Bai et al., 2023b), still employ the F1 score as the evaluation metric, which measures the similarity between the generated response of the large language model and the candidate answers. However, the generated content of large language models is highly diverse, and the candidate answers cannot cover all possible correct responses, making the evaluation of large language models less accurate.

To address these challenges, we propose a new long context benchmark called **Marathon** based on Loogle (Li et al., 2023b), which includes 6 evaluation tasks for assessing the ability of large language models. The context lengths in the benchmark range from 60k to 260k+. For each long context, there is a related question and four carefully crafted answer options. The large language model needs to provide the correct answer option based on the

<sup>1</sup>We also provide an online evaluation website: <https://openbenchmark.online/marathon>

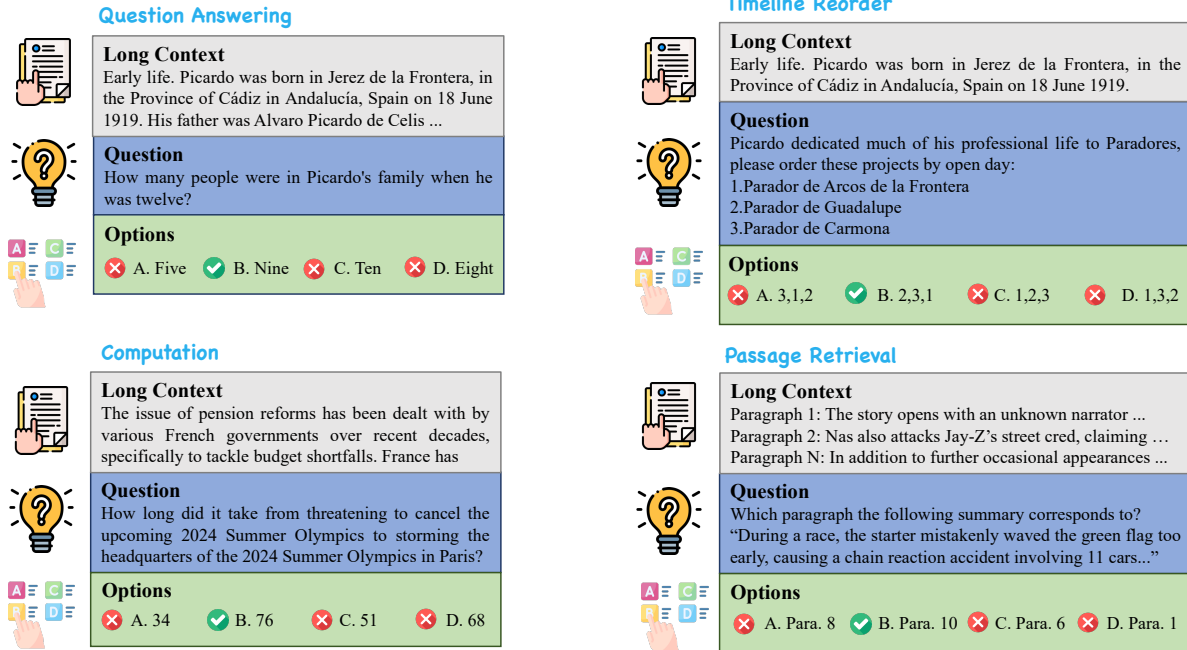


Figure 1: An example of a test case in benchmark. The context is truncated for display purposes.

given long context.

We conducted evaluations on **Marathon** using various large language models with different parameter sizes and context windows, combined with different methods of compressing long contexts. We found that the long context processing capabilities of most large language models still need improvement. In fact, the 7B Mistral (Jiang et al., 2023a) model outperformed the 70B Beluga (Mahan et al., 2023) model in terms of long context understanding. Some compression methods for long contexts were found to be effective in improving large language models' understanding of long contexts in certain cases.

## 2 Related Work

### 2.1 Flash Attention

The training and inference of ultra-long sequences have been hindering the further development of large language models. In traditional Attention mechanisms, the GPU memory consumption of Attention Scores is proportional to the square of the sequence length. Therefore, when the sequence length reaches a certain threshold, both training and inference will encounter out-of-memory (OOM) issues. The emergence of Flash Attention (Dao et al., 2022) allows us to avoid computing the Attention Scores matrix and directly compute the final Attention Output result, greatly reducing the GPU mem-

ory required during model training and inference. Moreover, thanks to the reduction in the number of memory accesses between HMB (high-bandwidth memory) and GPU SRAM (static random-access memory) brought by Flash Attention, the speed of model training and inference is accelerated. Flash Attention 2 (Dao, 2023) further optimizes the implementation of Flash Attention, further accelerating the speed of model computation. Currently, Flash Attention 2 (Dao, 2023) has become a key technology for optimizing the training and inference process of large language models, playing a crucial role in modeling and inferring ultra-long sequences.

### 2.2 Long Context Models

The ability of large language models for handling long contexts has become increasingly important. ChatGPT (OpenAI, 2023a) supports a window size of 16k, while GPT-4 supports a window size of 128k, and Claude-2.1 supports a window size of 200k<sup>2</sup>. Many open-source large language models have started to expand the size of their context window. Longchat (Li et al., 2023a) and MPT (Team, 2023) have achieved a window size of 16k, while Mistral (Jiang et al., 2023a) and Zephyr (Tunstall et al., 2023) have achieved a window size of 32k. By utilizing an adapted Rotary Embedding (Su et al., 2022) and sliding window (Beltagy et al.,

<sup>2</sup><https://www.anthropic.com/index/claude-2-1>

2020) during fine-tuning, MistralLite, based on Mistral, has achieved a window size of 128k, enabling large language models to handle even longer contextual information.

### 2.3 Prompt Compression

Although larger context windows enable large language models to handle longer contextual information, processing long-context information requires a significant amount of computing resources and places high demands on hardware. It also necessitates longer computational time, even in the inference stage. Therefore, some methods like LLM-Lingua (Jiang et al., 2023c) and LongLLMLingua (Jiang et al., 2023b) have been proposed to compress long contexts.

### 2.4 Retrieval Augmented Generation

Retrieval Augmented Generation (RAG) was originally proposed and applied to NLP tasks in (Lewis et al., 2020), and it has now become a mainstream method for improving the generation capability of large language models. RAG can extract the most relevant data from external knowledge bases and hand it over to the large language model for processing. This can alleviate the hallucination problem of large language models and enable people to trace the source of the content generated by large language models, ensuring the reliability of the generated content. Additionally, RAG can also be used to extract information from long documents that is most relevant to the user’s query. This ensures that key information required to provide correct answers to questions is not lost while reducing the length of the context. Many projects such as longchain (longchain, 2022) and LlamaIndex (Liu, 2022) have achieved significant progress in combining RAG with large language models, greatly facilitating related research in this direction.

### 2.5 Long Context Benchmarks

There have been many recent benchmarks used to assess the long context processing ability of large language models, such as Loogle (Li et al., 2023b) and LongBench (Bai et al., 2023b). Liu et al. (2023) on the other hand, noticed that the position of key information in long contexts greatly affects the capability of large language models to correctly understand and process text. Therefore, they used the NaturalQA (Kwiatkowski et al., 2019) dataset to construct a new benchmark to test the impact of different positions of key information in long

Task	No. Samples
Comprehension and Reasoning	357
Multiple Information Retrieval	341
Timeline Reorder	152
Computation	97
Passage Retrieval	300
Short Dependency QA	283
<b>Total</b>	<b>1530</b>

Table 1: Statistics of Marathon.

context on the text processing capability of large language models.

Although Loogle (Li et al., 2023b) and LongBench (Bai et al., 2023b) have constructed a relatively comprehensive set of evaluation tasks, the evaluation metrics used are still F1-score, Bleu or Rouge, which cannot accurately evaluate the ability of large language models to handle and understand long contexts.

## 3 Marathon

Currently, the mainstream benchmarks for evaluating large language models mostly adopt a multiple-choice question format, such as MMLU (Hendrycks et al., 2021) and C-Eval (Huang et al., 2023). The use of multiple-choice questions can effectively avoid scenarios where large language models generate correct answers but receive low scores due to the lack of relevant representations in the reference answers, or generate correct answers but receive high scores due to their similarity with the content of the reference answers. Therefore, based on Loogle (Li et al., 2023b) and Longbench (Bai et al., 2023b), we have constructed a multiple-choice long-text benchmark to more reasonably assess the long-context comprehension ability of large language models.

### 3.1 Overview

The Marathon benchmark consists of 6 tasks: *Comprehension and Reasoning*, *Multiple Information Retrieval*, *Timeline Reorder*, *Computation*, *Passage Retrieval*, and *Short Dependency QA*. These tasks can be categorized into four types based on the nature of the questions: Question Answering, Timeline Reorder, Computation, and Passage Retrieval. Figure 1 shows sample test cases for each type of

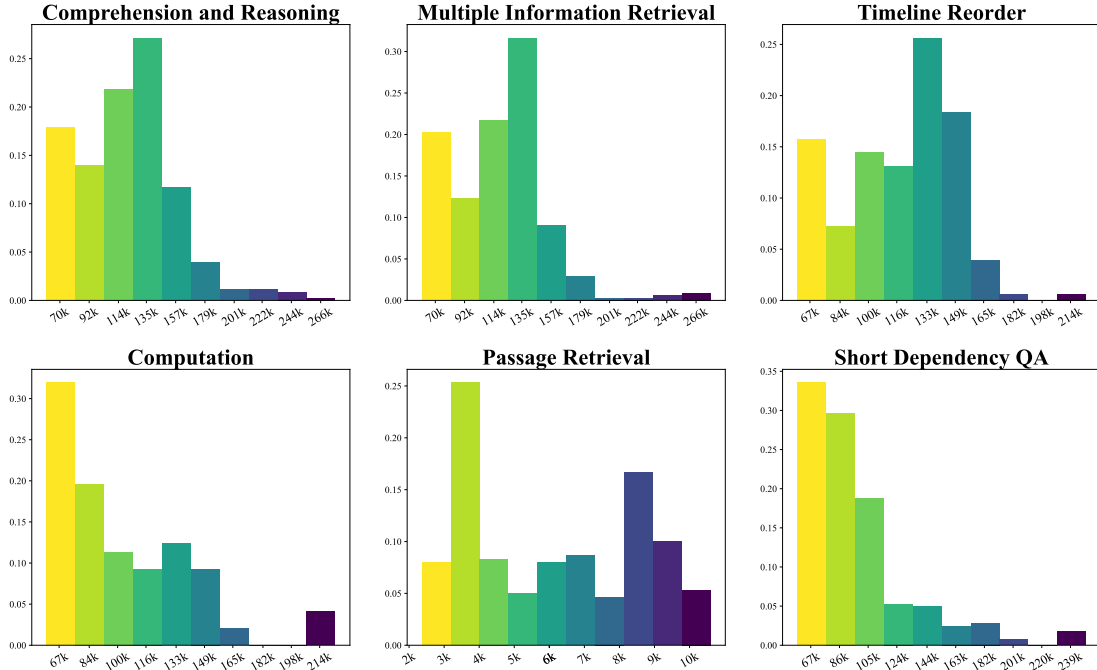


Figure 2: The distribution of context length in Marathon.

question. Table 1 shows the quantity information of test samples for each task.

### 3.2 Construction

All the test samples in the benchmark are in the form of multiple-choice questions, with each question containing one correct answer option and several distractor options. We use GPT-4 to generate the distractor options for each question. For each question, we divide the long context into multiple fragments of length 12,000 and randomly select one fragment. We require GPT-4 to generate three distractor options based on the given context fragment, question, and correct answer. The purpose of this approach is to avoid using excessively long context that exceeds GPT-4’s context window, which may affect the accuracy of the generated results. By using shorter contexts, we can obtain distractor options that are more relevant to these shorter contexts. Finally, to ensure the effectiveness and accuracy of these distractor options, we manually verify the options of each test sample.

### 3.3 Question Answering

*Comprehension and Reasoning*, *Multiple Information Retrieval* and *Short Dependency QA* are all types of traditional question-answer formats. The difference lies in the fact that *Comprehension and Reasoning*, *Multiple Information Retrieval* are se-

lected from the Long Dependency QA dataset in Loogle (Li et al., 2023b), while *Short Dependency QA* is selected from the Short Dependency QA dataset in Loogle (Li et al., 2023b). In the question answering tasks, each question is accompanied by a corresponding long context, and the large language model is required to infer the correct answer according to the long context. For the *Short Dependency QA* task, the relevant content for the correct answer is relatively concentrated within the long context. For *Comprehension and Reasoning* and *Multiple Information Retrieval* tasks, the content relevant to the correct answer is more scattered throughout the long context. Therefore, the large language model needs to possess strong long context understanding capability in order to solve the question correctly.

In the upper left of Figure 1, an example of a Question Answering task is provided. The question asks the large language model to answer a related question based on the content in the long context.

### 3.4 Timeline Reorder

*Timeline Reorder* task is a relatively novel question answering task. Unlike traditional question answering tasks, in the *Timeline Reorder* task, the question format requires large language models to sort a series of events described in a long context according to their chronological order. This

task aims to examine the large language models’ understanding of temporal relationships. Due to the dispersed distribution of events that need to be sorted by chronological order in the long context, large language models not only need to possess a correct understanding of temporal order but also require strong long context processing capabilities to answer correctly, which makes it a challenging task.

In the upper right of Figure 1, an example of the *Timeline Reorder* task is provided. The question requires the large language model to sort three events mentioned in the long context according to their chronological order.

### 3.5 Computation

*Computation* task is also different from traditional question answering tasks. Its question format involves providing a question related to numerical computation and requires the large language model to perform numerical calculations based on relevant content in the long context. For example, it may require calculating the number of children a certain character has at a specific time point, considering that the long context describes the character’s life events, including the death of a child due to illness, which may affect the number of the character’s offspring at subsequent time points. Therefore, to answer this question correctly, the large language model not only needs to be able to perform ordinary numerical calculations but also needs to capture all the key information related to the question. Compared to traditional computation and question answering tasks, this task is more challenging and can better reflect the large language model’s capability to comprehend long context.

In the bottom left of Figure 1, an example of a *Computation* task is provided. The question requires the large language model to complete a numerical calculation question based on the content in the long context.

### 3.6 Passage Retrieval

*Passage Retrieval* task is one form of task in the Longbench (Bai et al., 2023b). In order to enhance the diversity of our benchmark tasks, we have sampled 300 test data from the Passage Retrieval task in Longbench (Bai et al., 2023b), and reformed them into multiple-choice format using the method mentioned above. We have incorporated this task into our benchmark. The *Passage Retrieval* task requires large language models to locate the para-

graph in a long context that corresponds to the given description in the question. Since the test data of the *Passage Retrieval* task is sampled from Longbench (Bai et al., 2023b), there are some limitations in terms of context length and timeliness. However, it remains a highly valuable task format. In future work, we will update its content to make it more suitable for the current needs of evaluating large language models.

A sample of the *Passage Retrieval* task is provided in the bottom right of Figure 1. The task requires large language models to locate the paragraph in a long context that corresponds to the given description in the question.

## 4 Experiments

### 4.1 Setup

**Models.** In this evaluation, we selected models with varying parameter sizes, ranging from 7B to 70B, and different context window sizes, ranging from 8k to 200k. The models evaluated in this study include ChatGLM3-6B-32K (Zeng et al., 2022; Du et al., 2022), Mistral-7B-Instruct-v0.1 (Jiang et al., 2023a), Zephyr-7B- $\beta$  (Tunstall et al., 2023), Longchat-13B-16K (Li et al., 2023a), Qwen-14B-Chat (Bai et al., 2023a), Yi-34B (01.AI, 2023), Alfred-40B-1023 (Hallström et al., 2023), StableBeluga-2-70B (Mahan et al., 2023), Tulu-2-DPO-70B (Iverson et al., 2023), ChatGPT-1106 (OpenAI, 2023a), and GPT-4-1106-preview (OpenAI, 2023b).

**Methods.** In this evaluation, we first assessed the inherent ability of various models to comprehend long contexts. Then, we evaluated the current mainstream methods for handling long contexts: Compression and RAG. Specifically, for the Compression method, we assessed LongLLMLingua (Jiang et al., 2023b), while for the RAG method, we evaluated two retrieval approaches using OpenAI Embedding and Jina Embedding respectively.

### 4.2 Implementation Details

**Prompt.** We used the same prompt template to ask questions for all models, and required the answers to be returned in JSON format. The specific prompt format can be seen in Figure 9.

**LongLLMLingua.** For LongLLMLingua, we set the *compression rate* to 0.5, the *dynamic context compression ratio* to 0.4. We also sort the compressed contexts based on their importance.

Model	Para.	CW	C&R	MIR	TR	Com.	PR	SDQA	Avg.
GPT-4	-	128K	77.03%	69.21%	69.08%	60.82%	100.00%	95.41%	78.59%
ChatGPT	-	16K	62.18%	51.32%	19.74%	34.02%	95.67%	81.27%	57.37%
<b>Vanilla</b>									
Chatglm	6B	32K	55.46%	46.63%	30.26%	37.11%	81.33%	79.51%	55.05%
Mistral	7B	32K	46.22%	41.94%	28.95%	23.71%	49.67%	48.41%	39.81%
Zephyr	7B	32K	41.46%	37.83%	33.24%	21.65%	47.67%	47.00%	37.97%
Longchat	13B	16K	37.25%	34.60%	28.29%	27.84%	42.00%	45.23%	35.87%
Qwen	14B	8K	45.38%	39.00%	26.32%	23.71%	56.33%	44.88%	39.27%
Yi	34B	200K	59.66%	47.21%	37.50%	36.08%	90.00%	65.02%	55.91%
Alfred	40B	8K	40.90%	39.30%	26.32%	20.62%	49.00%	47.70%	37.31%
Beluga	70B	4K	55.74%	43.70%	36.84%	36.08%	65.33%	59.36%	49.51%
Tulu2	70B	8K	46.50%	35.48%	30.26%	22.68%	46.33%	46.29%	37.92%
<b>LongLLMLingua Compression</b>									
Chatglm	6B	32K	47.06%	37.54%	25.66%	22.68%	98.33%	56.18%	47.91%
Mistral	7B	32K	40.06%	31.38%	23.03%	27.84%	57.00%	42.76%	37.01%
Zephyr	7B	32K	30.81%	26.39%	23.68%	18.56%	54.00%	27.92%	30.23%
Longchat	13B	16K	37.82%	28.74%	26.32%	20.62%	61.67%	38.52%	35.61%
Qwen	14B	8K	42.58%	36.66%	27.63%	26.80%	88.67%	42.40%	44.12%
Yi	34B	200K	49.58%	42.23%	30.26%	22.68%	90.33%	56.89%	48.66%
Alfred	40B	8K	38.94%	32.84%	26.32%	29.90%	59.00%	45.94%	38.82%
Beluga	70B	4K	50.42%	42.82%	36.84%	27.84%	94.00%	63.60%	52.59%
Tulu2	70B	8K	45.94%	35.19%	34.87%	12.37%	98.00%	53.00%	46.56%
<b>OpenAI Embedding Retrieval</b>									
Chatglm3	6B	32K	56.58%	43.40%	28.95%	28.87%	81.33%	66.78%	50.99%
Mistral	7B	32K	51.54%	47.21%	27.63%	27.84%	79.00%	67.84%	50.18%
Zephyr	7B	32K	52.38%	43.99%	28.29%	24.74%	76.67%	71.73%	49.63%
Longchat	13B	16K	38.10%	25.81%	19.08%	12.37%	43.00%	41.34%	29.95%
Qwen	14B	8K	61.34%	46.33%	31.58%	18.56%	93.00%	69.96%	53.46%
Yi	34B	200K	66.39%	55.13%	38.82%	42.27%	95.00%	83.75%	63.56%
Alfred	40B	8K	52.38%	48.39%	25.00%	27.84%	87.33%	67.14%	51.35%
Beluga	70B	4K	61.90%	46.33%	3.28%	21.65%	81.00%	75.27%	48.24%
Tulu2	70B	8K	64.99%	53.37%	41.45%	34.02%	95.67%	82.33%	61.97%
<b>Jina Embedding Retrieval</b>									
Chatglm	6B	32K	52.94%	44.57%	27.63%	23.71%	83.33%	71.38%	50.60%
Mistral	7B	32K	54.90%	43.99%	32.24%	25.75%	79.00%	76.33%	52.04%
Zephyr	7B	32K	52.66%	46.33%	30.92%	23.71%	91.00%	78.09%	53.79%
Longchat	13B	16K	42.58%	33.43%	22.37%	13.40%	57.67%	57.24%	37.78%
Qwen	14B	8K	60.50%	46.63%	44.08%	24.74%	94.33%	78.45%	58.12%
Yi	34B	200K	66.67%	54.25%	45.39%	38.14%	95.00%	83.39%	63.81%
Alfred	40B	8K	50.42%	44.28%	27.63%	25.77%	88.33%	71.02%	51.24%
Beluga	70B	4K	59.94%	49.85%	23.68%	27.84%	96.00%	77.03%	55.72%
Tulu2	70B	8K	64.99%	54.25%	38.82%	31.96%	95.00%	84.10%	61.52%

Table 2: The evaluation results of models on Marathon benchmark. C&R refers to *Comprehension and Reasoning* task; MIR refers to *Multiple Information Retrieval* task; TR refers to *Timeline Reorder* task; Com. refers to *Computation* task; PR refers to *Passage Retrieval* task; SDQA refers to *Short Dependency Question Answering* task; Avg. denotes the average accuracy across all tasks.

Type	Chatglm	Mistral	Zephyr	Longchat	Qwen	Yi	Alfred	Beluga	Tulu2
<b>Vanilla</b>									
<b>JSON</b>	69.48%	77.22%	84.51%	92.29%	62.29%	38.95%	13.27%	22.48%	30.72%
<b>JSON-like</b>	30.52%	21.18%	6.86%	3.99%	0.72%	28.56%	81.96%	0.33%	46.47%
<b>Plain Text</b>	0.00%	6.60%	8.63%	3.73%	36.99%	32.48%	4.77%	71.19%	22.81%
<b>LongLLMLingua Compression</b>									
<b>JSON</b>	94.58%	68.10%	63.20%	93.92%	90.92%	48.10%	13.66%	29.97%	35.45%
<b>JSON-like</b>	5.42%	22.68%	9.15%	2.81%	0.06%	26.60%	85.95%	0.59%	43.46%
<b>Plain Text</b>	0.00%	9.22%	27.65%	32.68%	9.02%	25.29%	0.39%	69.54%	19.08%
<b>OpenAI Embedding Retrieval</b>									
<b>JSON</b>	52.88%	84.31%	21.83%	31.11%	65.62%	67.71%	16.27%	6.27%	98.43%
<b>JSON-like</b>	42.42%	6.67%	69.87%	23.73%	16.93%	32.16%	83.73%	0.00%	1.11%
<b>Plain Text</b>	4.71%	9.02%	8.30%	45.16%	17.45%	0.13%	0.00%	93.73%	0.46%
<b>Jina Embedding Retrieval</b>									
<b>JSON</b>	86.34%	83.14%	17.91%	32.75%	63.33%	65.69%	0.00%	8.43%	97.19%
<b>JSON-like</b>	9.15%	6.21%	73.86%	21.70%	17.91%	34.18%	100.00%	0.007%	2.16%
<b>Plain Text</b>	4.51%	10.65%	8.24%	45.56%	18.76%	0.13%	0.00%	91.50%	0.65%

Table 3: The evaluation results of large language models on the Marathon benchmark for instruction following.

**Embedding Retrieval.** For Embedding Retrieval, we utilize the ServiceContext and VectorStoreIndex of the Llama-Index (Liu, 2022). We employ various models as LLMs (Language Models), testing the OpenAI Embedding model and the Jina Embedding model as Embedding Models respectively. The default parameter settings are retained, with a chunk size of 1024 and a top-k value of 2. As for Jina Embedding, we set the pooling method to "mean" to align with Jina’s encode implementation.

**Hardware.** All experiments in this evaluation were conducted on a server with 4\*A100 80GB.

## 4.3 Results

### 4.3.1 Main results

Table 2 presents the performance of different models with different optimization methods on various tasks in the evaluation. For a more intuitive comparison of the results, one can refer to Figure 3, Figure 4, Figure 5, Figure 6, Figure 7, and Figure 8 in appendix. It can be observed that both OpenAI Embedding Retrieval and Jina Embedding Retrieval outperform LongLLMLingua.

In addition, all models have lower accuracy in the *Timeline Reorder* task and *Computation* task compared to other tasks. The adoption of the LongLLMLingua method did not lead to improvement, and the improvement achieved by the opti-

mization method RAG was also limited.

**Vanilla.** Among the Vanilla methods, the highest accuracy is achieved by Yi-34B, which has 34B parameters but achieves an accuracy of 55.91%. Followed by ChatGLM3-6B-32K, which achieves an accuracy of 55.05% despite having only 6B parameters. Next is Beluga-70B, which achieves an accuracy of 49.51% despite having a context window of only 4K. The average accuracy of the remaining models does not differ significantly, and none of them exceeds 40%.

**LongLLMLingua.** Compared to the Vanilla method, LongLLMLingua only improved the accuracy of Qwen by 4.85%, Alfred by 1.51%, Beluga by 3.08% and Tulu2 by 8.64%. However, it reduced the accuracy of ChatGLM3 by 7.14%, Mistral by 2.8%, Zephyr by 7.74%, Longchat by 0.26%, and Yi by 7.25%.

**OpenAI Embedding Retrieval.** Compared to the Vanilla method, OpenAI Embedding Retrieval improves the accuracy of Mistral by 10.37%, Zephyr by 11.66%, Qwen by 14.19%, Yi by 7.65%, Alfred by 14.05% and Tulu2 by 24.05%. However, it results in a decrease in the accuracy of ChatGLM3 by 4.06%, Longchat by 5.92% and Beluga by 1.27%.

**Jina Embedding Retrieval.** Compared to the Vanilla method, Jina Embedding Retrieval has improved the accuracy of most models. Specifically, Mistral’s accuracy has improved by 12.23%, Zephyr’s accuracy has improved by 15.82%, Longchat’s accuracy has improved by 1.91%, Qwen’s accuracy has improved by 18.15%, Yi’s accuracy has improved by 7.9%, Alfred’s accuracy has improved by 13.93%, Beluga’s accuracy has improved by 6.21%, and Tulu2’s accuracy has improved by 23.60%.

### 4.3.2 Instruction Following Capability

During the evaluation process, we found that many models had poor instruction-following capabilities. We specified in the prompt that the large language model should respond in JSON format and provided an example of the desired response. However, sometimes the model still responded in other formats or attempted to respond in JSON format but with incomplete results. We conducted a statistical analysis of the evaluation results of various models to estimate their instruction following capabilities, the results are displayed in table 3.

JSON represents the correct generation of a JSON format response by the large language model. JSON-like refers to the attempt by the large language model to generate a response in JSON format but, for various reasons (such as exceeding length limits and being truncated, or having formatting errors), it fails to generate a correct JSON format response. Plain Text indicates that the large language model uses a different format for the response.

From Table 3, we can see that even Yi, which achieved the highest answer accuracy rate, only had an instruction-following capability of 38.95%. Beluga, on the other hand, had an instruction-following capability of only 22.48%. However, Longchat, despite having a low question-answering accuracy rate, exhibited strong instruction-following capabilities with a score of 92.29%, second only to ChatGPT’s 99.61%. Mistral and Zephyr also showed instruction-following capabilities of 77.22% and 84.51%, respectively.

## 5 Discussion

**JSON Format.** At the recent OpenAI Developer Day<sup>3</sup>, OpenAI announced two major updates for GPT-4. One is parallel function invocation and the

<sup>3</sup><https://devday.openai.com>

other is the ability to specify response format, requiring GPT-4 to return responses in JSON format. These two features are fundamental for building AGI (Artificial General Intelligence). With parallel function invocation, large language models can simultaneously call multiple utility functions to complete user tasks. The ability to respond in JSON format ensures the smooth inclusion of parameters and result returns for function invocation.

**Document Context.** After the OpenAI Developer Day, GPT-4 acquired a new feature called Knowledge Retrieval, which allows users to upload their own files and retrieve relevant information from these files to respond to user queries, which is an important application scenario for RAG.

In the future, more and more large language models will implement similar functionalities. The evaluation format for long context question answering tasks will also change. Instead of directly inputting the long context in the prompt, the benchmarking work for long-context question-answering should be aligned with the direction of answering user questions based on the files provided by the user.

**Online Benchmarks.** With the continuous development of large language models, the methods and benchmarks for evaluating them also need to be changed accordingly. In the past, when evaluating large language models, we relied on a research team or individual to propose a static benchmark, and then conducted evaluations based on this benchmark. However, the data in static benchmarks would leak over time and be used for model training, resulting in less accurate evaluation results. Moreover, it is inefficient for a research team or individual to build benchmarks, and it is also difficult for researchers to continuously find new evaluation data.

The power of an individual is limited, but the power of a collective is infinite. We call for the construction of Online Benchmarks, which will gather the strength of the entire open-source community to collectively maintain a continuously updated benchmark. This benchmark will have regularly updated evaluation tasks and periodically replaced evaluation formats to cope with the rapid development of large language models.

## 6 Conclusion

In this experiment, we compared six open source large language models with different parameter



sizes and different context windows, as well as OpenAI’s ChatGPT and GPT-4. We also compared the optimization methods based on Compression, LongLLMLingua, and the optimization methods based on RAG, OpenAI Embedding Retrieval, and Jina Embedding Retrieval.

The experimental results show that the method based on RAG is more effective than the method based on Prompt Compression, and Jina Embedding performs better than OpenAI Embedding. Furthermore, the instruction-following ability of open-source large language models is generally poor, despite their ability to achieve high accuracy in question-answering tasks.

## References

- 01.AI. 2023. [Yi-series](#).
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023a. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- 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. 2023b. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. [Longformer: The long-document transformer](#).
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri 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, Josh 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](#).
- Tri Dao. 2023. FlashAttention-2: Faster attention with better parallelism and work partitioning.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems*.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335.
- Oskar Hallström, Amélie Chatelain, Clément Thiriet, Julien Séailles, Adrien Cavallès, and Axel Marmet. 2023. [Alfred-40b-1023](#).
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#).
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023. [C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models](#).
- Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023. [Camels in a changing climate: Enhancing lm adaptation with tulu 2](#).
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023a. [Mistral 7b](#).
- Huiqiang Jiang, Qianhui Wu, , Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023b. [Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression](#). *ArXiv preprint*, abs/2310.06839.
- Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023c. [Lmlingua: Compressing prompts for accelerated inference of large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS'20*, Red Hook, NY, USA. Curran Associates Inc.
- Dacheng Li, Rulin Shao, Anze Xie, Ying Sheng, Lianmin Zheng, Joseph E. Gonzalez, Ion Stoica, Xuezhe Ma, and Hao Zhang. 2023a. [How long can open-source llms truly promise on context length?](#)
- Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. 2023b. [Can long-context language models understand long contexts?](#)
- Jerry Liu. 2022. [LlamaIndex](#).
- 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. [ArXiv:2307.03172](#).
- longchain-ai longchain. 2022. [LangChain](#).
- Dakota Mahan, Ryan Carlow, Louis Castricato, Nathan Cooper, and Christian Laforte. 2023. [Stable beluga models](#).
- Ha-Thanh Nguyen. 2023. [A brief report on lawgpt 1.0: A virtual legal assistant based on gpt-3](#).
- OpenAI. 2023a. [Chatgpt](#). <https://chat.openai.com>.
- OpenAI. 2023b. [Gpt-4 technical report](#).
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2022. [Roformer: Enhanced transformer with rotary position embedding](#).
- MosaicML NLP Team. 2023. [Introducing mpt-7b: A new standard for open-source, commercially usable llms](#). Accessed: 2023-05-05.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. [Zephyr: Direct distillation of lm alignment](#).
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. [Glm-130b: An open bilingual pre-trained model](#). *arXiv preprint arXiv:2210.02414*.

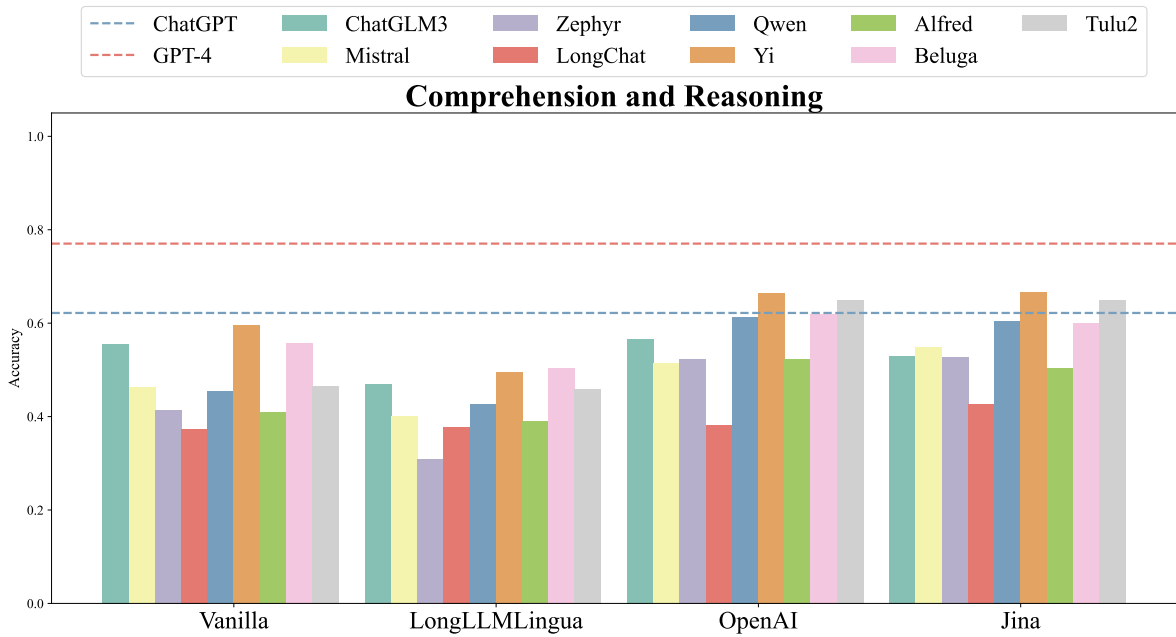


Figure 3: The performance of models on comprehension and reasoning task.

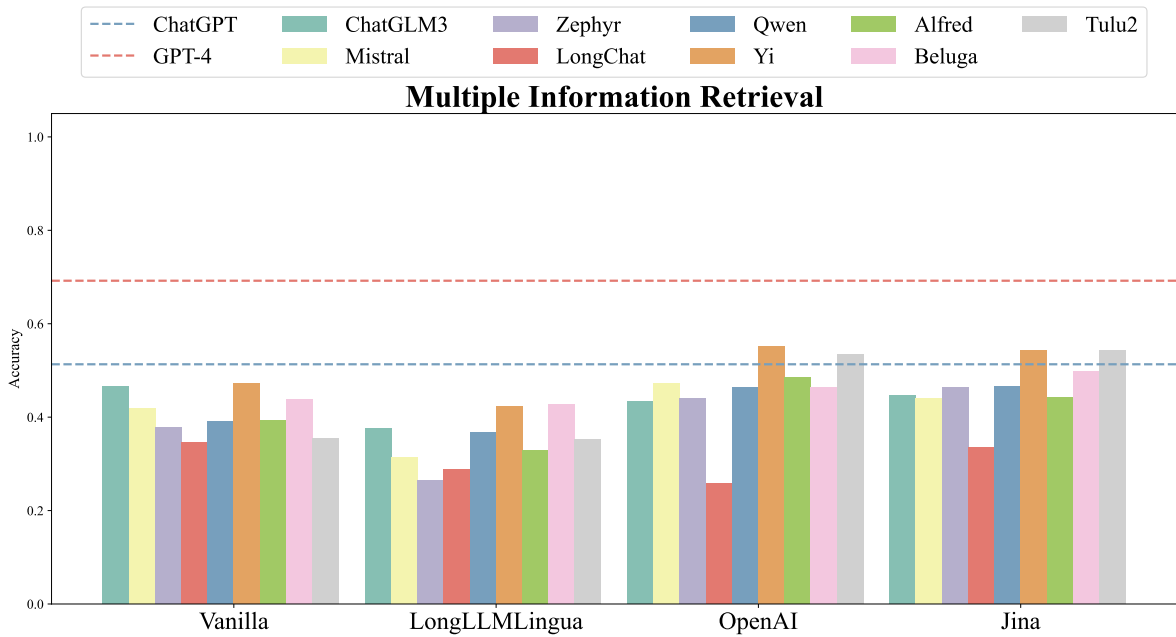


Figure 4: The performance of models on multiple information retrieval task.

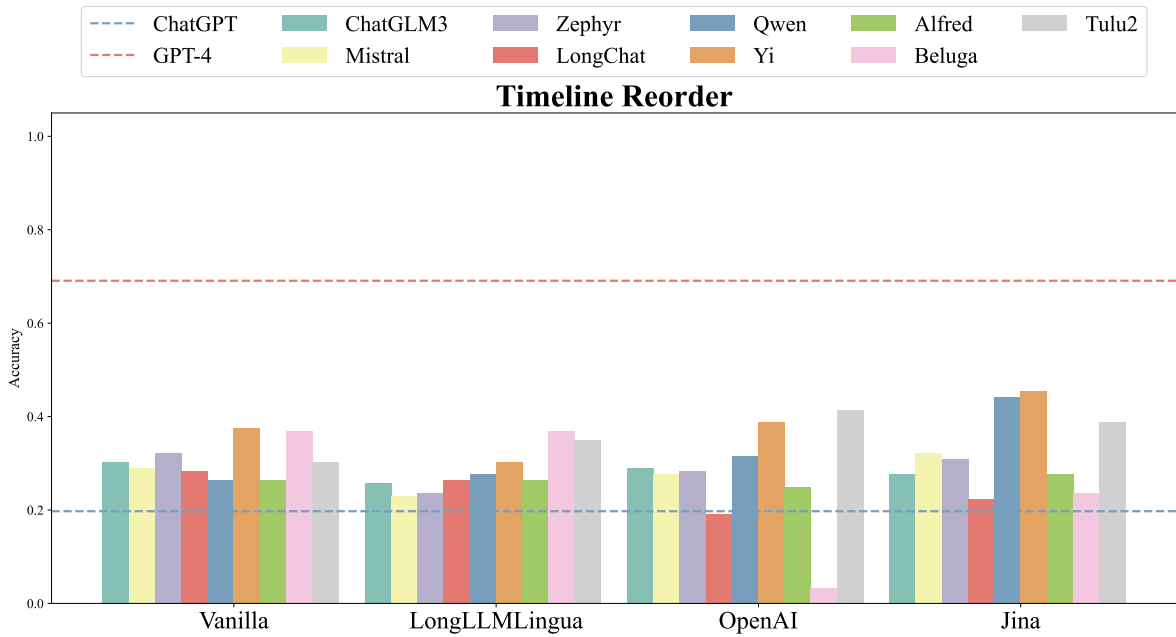


Figure 5: The performance of models on timeline reorder task.

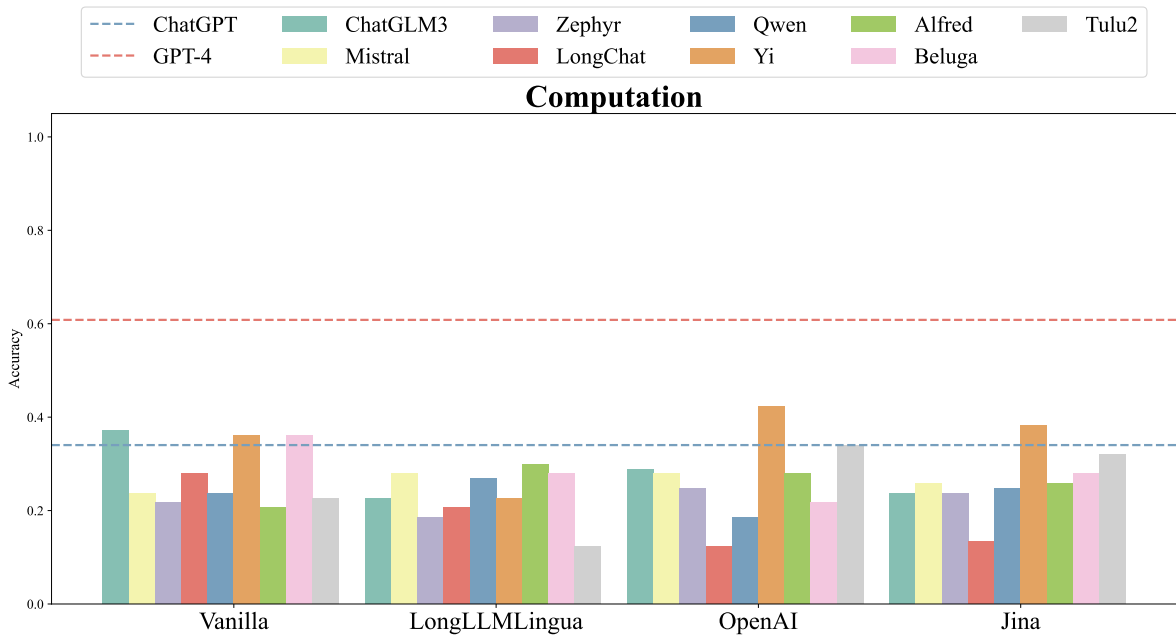


Figure 6: The performance of models on computation task.

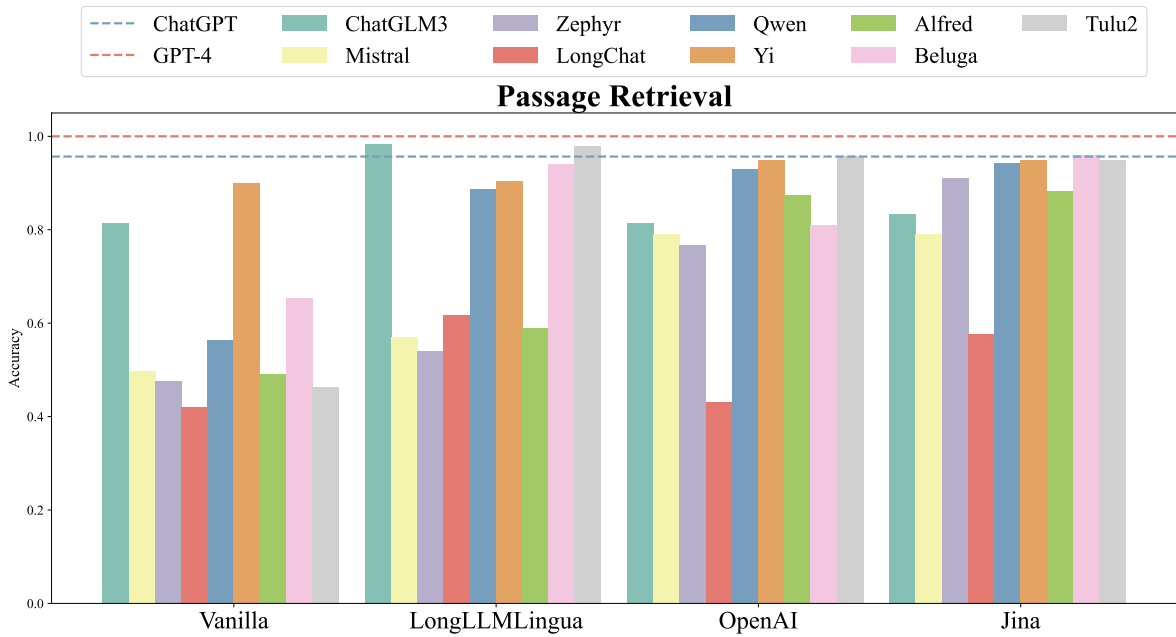


Figure 7: The performance of models on passage retrieval task.

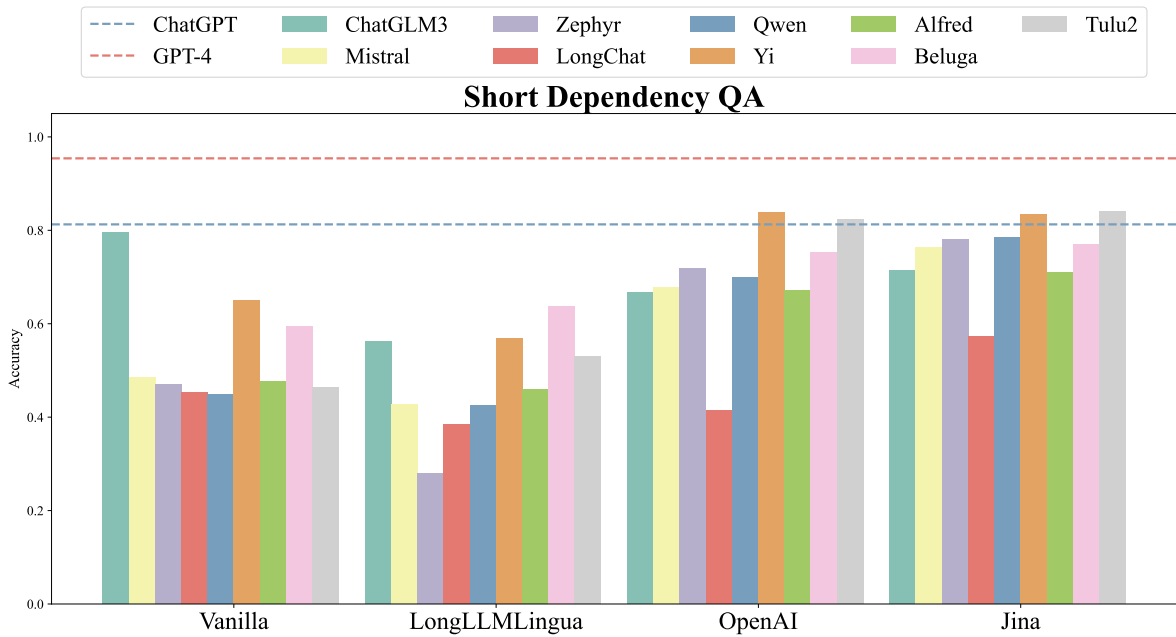


Figure 8: The performance of models on short dependency question answering task.

You are an expert at reading and analyzing lengthy texts for examinations. Your task is to carefully read the provided text, understand its content and details, and accurately answer multiple-choice questions about the text. Keep in mind that the correct answer must be based entirely on the content of the text, without including any external information or personal opinions.

**Context:**

Olympics on the Whistler Sliding Centre in Whistler, British Columbia, Canada. Hours later, the International Luge Federation concluded that the accident was caused by a steering error and not a track error; nevertheless, ...

**Question:**

Based on the description above, what is the name of son of lord krishna?

**Options:**

- A. Jon Owen
- B. Nodar Kumaritashvili
- C. Ulrich Hahn
- D. Paul Aste

Please answer this question with JSON format, for example {"option": "A"}.

**Answer:**

Figure 9: An example of a prompt. The context is truncated for display purposes.