

# Attention-Guided Context Pruning for Large Language Models

Anonymous ACL submission

## Abstract

While large language models (LLMs) demonstrate remarkable capabilities, their effectiveness is hindered by the ever-increasing length of prompts, which introduces information sparsity and substantial computational overhead. Existing prompt pruning methods, such as LLMLingua, lack contextual awareness and offer limited flexibility in controlling compression rates, often resulting in either insufficient pruning or excessive information loss. In this paper, we propose AttentionPrompt, an attention-guided prompt pruning method for LLMs in the RAG setting. The core idea of AttentionPrompt lies in its attention focus mechanism, which reformulates user queries into a next-token prediction paradigm. This mechanism isolates the query’s semantic focus to a single token, enabling precise and efficient attention calculation between queries and contexts. Extensive experiments on LongBench and Babilong benchmarks show that AttentionPrompt achieves up to 6.3x context compression while outperforming LLMLingua methods by around 10% in key metrics.

## 1 Introduction

Large language models demonstrated remarkable capabilities in a variety of applications such as reasoning (Huang and Chang, 2023), question-answering (Rajpurkar et al., 2016; Wang et al., 2024), and open-ended text generation (Que et al., 2024). While effective, LLMs encounter significant challenges in handling long contexts. As the context becomes excessively long, it brings large amounts of redundant and irrelevant information. This issue, as highlighted by Shi et al. (2024), can lead to hallucinations and a decline in the performance of the LLM (Chiang and Cholak, 2022).

To mitigate these issues, recent work has explored context compression techniques (Jiang et al., 2023; Cheng et al., 2024; Verma, 2024). For example, LLMLingua (Jiang et al., 2023) compresses

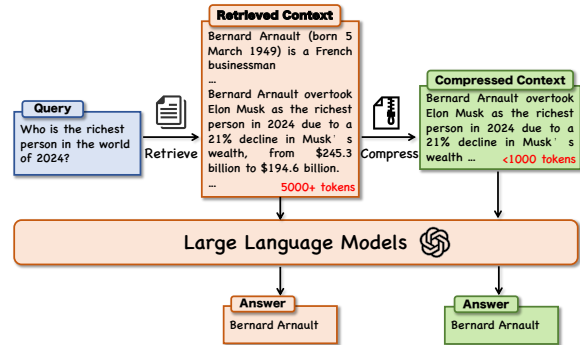


Figure 1: Illustration of LLM context compression in the RAG setting

prompts by using a budget controller that allocates varying compression ratios to different components of the prompt, such as instructions and demonstrations. However, these methods lack context-awareness, making it challenging to determine the optimal compression ratio for a given LLM, resulting in context redundancy or over-compression.

Although attention mechanisms are central to LLMs for content selection, there is limited research on using them specifically for context pruning. This gap is primarily due to two fundamental challenges that hinder the effectiveness of attention in the inference pipeline: (1) The long nature of contexts exacerbates attention dilution (Hsieh et al., 2024; Zhu et al., 2024). When the context is excessively long, the attention scores are spread thin, causing the model to allocate less focus to any single token. This dilution of attention reduces the ability of the model to concentrate on the most relevant parts of the context; and (2) The sentence-based nature of the query makes it challenging to directly align with critical content in the context. The attention score distribution varies across tokens within a sentence, with some tokens focusing on semantic information while others attend to details (Clark et al., 2019; Gu et al., 2022). For instance, in the sequence “Mary is in the car,” the

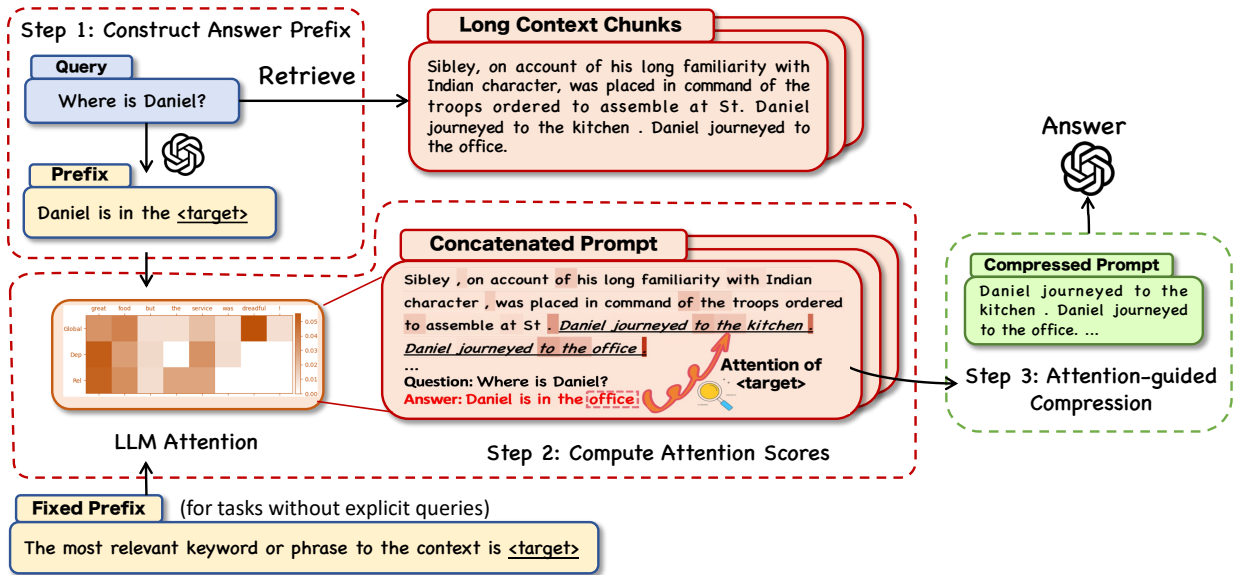


Figure 2: Illustration of AttentionPrompt in the typical RAG setting: 1) We first generate a hint prefix using the original query; if generation fails, we fallback to a fixed hint prefix. 2) The long context is then divided into smaller chunks for localized processing. 3) For each chunk, we batch the input with the hint prefix (1 token only) and use an LLM to compute attention features over the full context, focusing on the generated anchor token. 4) Using these attention features, we identify and extract the most relevant content from each chunk to construct a compressed context.

attention score of “in” tends to focus on semantic relationships, while “car” focuses more on a specific detail. Therefore, the query lacks the explicit or semi-structured reference points (such as keywords or pivotal terms) needed to guide the attention mechanism toward the most salient parts of the context.

In this paper, we propose a novel methodology called AttentionPrompt (Attention-Guided Prompt Pruning), which extracts the most salient information from the context by leveraging attention scores from intermediate layers of LLMs. To enable attention across queries and context, AttentionPrompt introduces an attention focus mechanism that isolates the query’s semantic focus to a single token, enabling precise and efficient attention computation between queries and contexts in a single pass. More specifically, for each query (e.g., “Where is Daniel?”), we construct an answer hint prefix (e.g., “Daniel is in the \_\_\_\_”) in a <next-token-prediction> format. For tasks (like summarization) that do not contain an explicit query. We employ a fixed prefix: “The most relevant keyword or phrase to the context is \_\_\_\_”. Next, we take the context, the query, and the constructed answer prefix as input to an LLM. The missing next token in the prefix stands for the focal point of the query, guiding the LLM’s attention to each token in the context. Finally, we produce a compressed context

by selecting sentences from the original context with the top-k attended tokens.

Extensive experiments on LongBench (Bai et al., 2024) and BABILong (Kuratov et al., 2024) demonstrate that our method achieves up to 6.3x context compression while maintaining or exceeding the performance of uncompressed contexts. The results suggest that AttentionPrompt not only facilitates the extraction of relevant information but also enhances the model’s reasoning capabilities. Particularly, it achieves these benefits without requiring additional training, making it highly adaptable across different models and practical for real-world applications.

Our key contributions are as follows:

- We propose a lightweight, transferable, and query-aware method for long-context pruning.
- We introduce a novel attention focus mechanism by reformulating user queries into a next-token prediction template, enabling precise and efficient computations of attention between queries and LLM contexts.
- We conduct extensive experiments on LongBench and Babilong benchmarks. Results demonstrate the effectiveness of AttentionPrompt in long-context LLM applications.

## 2 Preliminaries

**Attention Mechanism in LLMs** The attention mechanism is a key component in modern LLMs, allowing the model to focus on different parts of the input sequence (Vaswani et al., 2023). For a given context  $c$ , the mechanism calculates a score for each token  $t \in c$  based on its relevance to other tokens:

$$\text{Attention}_l(c, t) = \text{Softmax} \left( \frac{Q_l K_l^T}{\sqrt{d_k}} \right) V_l \quad (1)$$

where:  $Q_l$  is the query matrix at layer  $l$ ,  $K_l$  is the key matrix at layer  $l$ ,  $V_l$  is the value matrix at layer  $l$ ,  $d_k$  is the dimensionality of the key vectors.

The scores indicate how much “focus” the token receives from the model, providing insights into which tokens are most relevant in a given context. This enables it to be used for text compression, where selecting tokens with high attention scores can reduce the input to the model while retaining the most important information (Tarzanagh et al., 2023). Although attention is central to LLMs, it suffers from dilution in handling long contexts (Hsieh et al., 2024).

## 3 Method

In this work, we propose a novel approach to compressing the LLM context without compromising performance.

**Problem Formulation** Consider a typical RAG-based LLM inference scenario. Given a query  $q$  and a long context (including instructions, retrieved documents)  $c$ , a large language model  $P_\phi(a|q, c)$  takes the query and the context  $c$  as input and produces an answer  $a$  to the query. Our goal is to compress the long context  $c$  into a dense one  $c'$  such that  $|c'| \ll |c|$  while the LLM maintains the quality of responded answers when taking  $c'$  as input.

In this paper, we propose a novel attention-guided context pruning method called Attention-Prompt. The key idea of AttentionPrompt is reformulating each query into a next-token-prediction template (called answer hint prefix). This strategy allows the LLM to calculate the query-context attention through one token, therefore significantly improving the alignment between query and context, and reducing the time complexity for attention calculation.

Figure 2 shows the overall structure of our method. The pipeline involves three key steps: First, we generate an answer hint prefix for each query (Section Construct Answer Hint Prefix); Next, the generated prefix is appended to the original query and context as input to the LLM. The LLM is instructed to predict the follow-up token to the answer prefix, obtaining the attention scores (Section Compute Attention Features from LLM). Finally, we perform attention-guided compression: using attention scores from a language model, we identify and retain the most relevant parts of the long context. Each of the steps is elaborated in the following sections.

### 3.1 Construct Answer Hint Prefix

To improve alignment between the query and its context, we associate each query with an answer hint prefix that enables the LLM to compute the query-context attention through a single focal token.

For each query, an *answer hint prefix* is defined as an incomplete answer to the query in a *next-token-prediction* format, where the blank token to be predicted serves as the focal token of the query, directing the LLM’s attention to the most relevant parts of the context. For instance, as shown in Figure 2, for the query “Where is Daniel?”, the corresponding answer hint prefix can be “Daniel is in the \_\_\_\_”. For tasks (like summarization) that do not involve an explicit query, we employ a fixed hint prefix: “The most relevant keyword or phrase to the context is \_\_\_\_” that prompts the model to predict the most salient information (e.g., keyword) from the context. This ensures the generalization of our method across various tasks. When we take the context, query, and the answer hint prefix as input to the LLM, the token to be predicted by the model, such as “park”, becomes the focal point of the query, directing the model’s attention to the most relevant parts of the context and ensuring that the attention calculation focuses on the crucial information relevant to the query.

We prompt the LLM to automatically construct the answer hint prefix<sup>1</sup>. In detail, according to the query’s grammatical attributes, we categorize the answer hint prefix into two types: empty and non-empty ones. For example, for queries like wh-questions, the hint prefix can be derived from the query itself. In contrast, queries like

<sup>1</sup>We use GPT-4o Mini for the generation, the prompt detail and fixed hint prefix can be referred to Appendix.

yes/no-questions, where the answer’s first token is "Yes/No", already align with the next-token-prediction paradigm. Leveraging the semantic understanding capabilities of LLMs, we prompt them with example answer hint prefixes of various types of queries to automatically determine the type and generate the answer prefix hint.

By focusing attention on the focal token, this approach enhances both the precision and efficiency of the attention mechanism. The single-token focus accelerates computations by reducing the number of tokens involved in attention calculations. Simultaneously, concentrating on a target token—such as "car" in the sentence "Daniel is in the car"—enables the model to effectively identify and prioritize the most relevant information in the context.

### 3.2 Compute Attention Features from LLM

In this part, we aim to compute attention features based on the focal token. To address the issue of diluted attention scores, we divide the context  $C$  into smaller chunks. We adopt a uniform chunking strategy, assuming each chunk consists of  $m$  tokens. Let  $c_1, \dots, c_n$  represent the resulting chunks, where  $n = \lceil |C|/m \rceil$ . For each chunk  $c_j$ , we concatenate the chunked context, query, and instruction (we instruct the LLM to generate "none" after the hint prefix if the chunk is irrelevant, which will be used in Section [Compress with Attention](#)) and the generated answer hint prefix and feed this into the LLM. The LLM is then instructed to perform next-token prediction and compute the attention scores.

We define  $a_j$  as the first token generated following the answer hint prefix in chunk  $j$ . The attention feature  $A_j$  for  $a_j$  is computed as the sum of attention scores over all layers of the model, focusing on context tokens:

$$A_j = \sum_{l=0}^L \text{Attention}_l(c_j, a_j) \quad (2)$$

where  $L$  is the total number of layers in the model, and  $\text{Attention}_l(c_j, a_1)$  is the attention score at layer  $l$  for the token  $a_1$  relating to the context chunk  $c_j$ . This score is computed by the self-attention mechanism at each layer, capturing both local and global dependencies in the input.

The total attention feature  $A_j$  reflects the model’s focus on the most relevant components of the input when generating the first token, and is the sum of the attention values across all layers. We choose

to sum across all layers for analysis because the attention distribution in each layer can vary depending on the task. For easier tasks, earlier layers may already capture sufficient information to generate the final answer, while for more difficult tasks, the model might rely on the later layers (Jin et al., 2025). Since the function of each layer can vary from task to task, focusing on a single layer or a subset of layers could introduce bias in the attended information. To mitigate this issue, we sum the attention across all layers, which helps reduce task-specific bias. The choice of attention layers will be further explored in the ablation study in Section [Ablation Study](#).

### 3.3 Compress with Attention

For each chunk, after generating the focal token, we first check whether this token is "none." If this is the case, we skip the chunk, as it is deemed irrelevant to the task. If the focal token is valid, we proceed by identifying the tokens in the context that have the highest attention features with respect to the focal token. These attention features represent how much each token in the context is relevant to the focal token, which serves as the focal point of the query.

Next, we select the top- $k$  tokens based on their attention features, as these tokens are considered the most relevant to the focal token. To ensure that the context used for generating the final response is both relevant and concise, we focus on the sentences that contain these top- $k$  tokens. By selecting these sentences, we retain the information most pertinent to the focal token. These selected sentences are then concatenated to form a compressed context  $c'_j$ .

$$c'_j = \text{Concat}(\{s \mid t_r \in \text{Top-}k(A_j) \text{ and } t_r \in s\}) \quad (3)$$

where  $s$  denotes a selected sentence.

### 3.4 Time Efficiency and Batch Generation

Since we employ a next-token prediction paradigm, only one focal token needs to be generated for each chunk. Furthermore, as each chunk is processed independently, batch generation can be used to accelerate the process. This approach results in high time efficiency. Moreover, we can use quantified model to further accelerate the compression process. We provide the pseudocode of our method in Algorithm 1 in Appendix.

<b>Query:</b>	Where was the wife of Francis I Rákóczi born?	<b>Answer: Ozalj</b>
<b>Original Context:</b>	Passage 1: Waldrada of Lotharingia Waldrada was the mistress, and later the wife, of Lothair II of Lotharingia. Biography Waldrada’s family origin is uncertain. The prolific 19th-century French writer Baron Ernouf suggested that Waldrada was of noble Gallo-Roman descent, sister of Thietgaud... (7003 tokens)	city of Gyulafehérvár, Transylvania. ✗
<b>Random:</b>	drada., of. th-century French Baron Er of-R sister ofga bishopther arch of, not any socialoli,... (1400 tokens)	The text does mention it. ✗
<b>LLMLingua2:</b>	Waldrada Lotharingia mistress Lothair II Gallo sister Thietgaud niece Gunther Vita Sancti related Eberhard II Etichonids 855 Lothar II married Teutberga 858 862 Nicholas 863Charles ...(632 tokens)	Munkács.✗
<b>LongLLMLingua:</b>	Passage:Waldrada theressairia. is The proific 1th French Baron Ern thatadaoman, sister of Th Trier, Gun of and have suggested of social though anatic.itactoli thatada Ehard II,edbourgichon . ... (920 tokens)	Hungary. ✗
<b>AttentionPrompt:</b>	... Life Early years and family <b>Ilona was born Ilona Zrínyi in Ozalj</b> ... She was the daughter of Petar Zrinski, Ban (vicero) of Croatia, the niece of both Miklós Zrínyi and Fran Krsto Frankopan and <b>the wife of Francis Rákóczi I</b> ... (273 tokens)	Ozalj ✓

Table 1: Examples of compression results by various methods

## 4 Experimental Setup

In this experiment, we evaluate the efficacy of AttentionPrompt in long context compression.

### 4.1 Datasets

Due to fluctuations in the experimental results, each experiment was conducted three times, and the final result is the average of these trials.

**Long Context Reasoning** For our experiments, we incorporate datasets TriviaQA (Joshi et al., 2017), HotpotQA (Yang et al., 2018), and 2WikiMQA (Ho et al., 2020) from Longbench (Bai et al., 2024): TriviaQA assesses fact retrieval over long contexts, HotpotQA emphasizes multi-hop reasoning across dispersed clues, and 2WikiMQA tests the model’s ability to synthesize information from multiple sources.

**Babilong** (Kuratov et al., 2024): The BABI-Long benchmark provides a comprehensive framework for evaluating an LLM’s ability to reason and retrieve over long contexts. To assess AttentionPrompt’s capabilities on relatively shorter and sparse content, we select test splits of 1kqa1, 2kqa1, and 4kqa1 tokens to demonstrate its performance across different circumstances.

**Long Context Summarization** To demonstrate the robustness of AttentionPrompt, we also evaluate it on GovReport (Huang et al., 2021) from Longbench. In this task, we employ fixed hint prefix, showing the stability of our method.

### 4.2 Metrics

We measure the effectiveness of compression using three metrics:

**Exact Match (EM):** Measures the percentage of predicted answers that exactly match the ground-truth answers.

**LLM-as-a-judge scores:** We leverage GPT-4 (et al., 2024) to assess the correctness of the model-generated answers. Specifically, we input the query, the model-generated answer, and the ground truth answer into GPT-4o, asking it to determine whether the provided answer is correct. The prompt used for the scoring can be found in Appendix.

**Compression Ratio (CR):** The ratio of tokens in the original context to the compressed context, defined as  $CR = \frac{\# \text{ tokens in original context}}{\# \text{ tokens in compressed context}}$ .

### 4.3 Baselines

We compare AttentionPrompt with state-of-the-art context compression baselines, including LLMLingua2 (Pan et al., 2024) and LongLLMLingua (Jiang et al., 2024), which represent query-unaware and query-aware approaches, respectively. This diverse selection underscores the robustness of AttentionPrompt. For fair comparison, we match or exceed their compression rates. We also include the uncompressed context (“Original Prompt”), a random compression baseline (“Random”), and a vanilla RAG baseline using BGE-M3 (Chen et al., 2024).

### 4.4 LLMs for Compression

We select two open-source LLMs for evaluation, Llama-3.1-8B-Instruct (Dubey et al., 2024) and

376 Qwen-2.5-7B-Instruct (Yang et al., 2025). We ex- 423  
377 periment with both 8B and 70B versions of Llama- 424  
378 3.1-Instruct as our compression model, denoted as 425  
379 AttentionPrompt (8B) and AttentionPrompt (70B) 426  
380 respectively. Detailed hyperparameter configura- 427  
381 tions are provided in the Appendix. 428

## 382 5 Results 429

### 383 5.1 Overall Results 430

384 As the results in Tables 2 and 3 show, Attention- 431  
385 Prompt outperforms nearly all baseline methods 432  
386 across the benchmarks, consistently achieving su- 433  
387 perior results. On the LongBench benchmark, 434  
388 which includes contexts of tens of thousands of 435  
389 tokens, AttentionPrompt outperforms the LLMLin- 436  
390 gua methods in most metrics, particularly in BA- 437  
391 BILong 1k, where it achieves an 18% higher score 438  
392 than the best LLMLingua model. Specifically, in 439  
393 TriviaQA, AttentionPrompt outperforms the un- 440  
394 compressed method by 2% in EM score, demon- 441  
395 strating the effectiveness of our approach. How- 442  
396 ever, the low compression ratio is due to two parallel 443  
397 factors: the scattering of relevant information across 444  
398 the context and the use of uniform hyperparameters. 445  
399 All these factors contribute to the lower ratio, which 446  
400 we aim to enhance. We discuss this in more detail 447  
401 in the Ablation Study. Also, our method excels in 448  
402 summarization tasks, showing the stability of our 449  
403 method. Similar results can be observed in the BA- 450  
404 BILong benchmark, which involves a wider range 451  
405 of context lengths. Our method consistently out- 452  
406 performs all LLMLingua methods while achieving 453  
407 the highest compression ratio, further demonstrat- 454  
408 ing the effectiveness of AttentionPrompt in varying 455  
409 context sizes. 456

### 410 5.2 Efficiency Analysis 457

411 As discussed in Section Method, we divide the con- 458  
412 text into smaller chunks and use attention mech- 459  
413 anisms for compression. The overall time cost is 460  
414 spent in the forward pass of the compression model 461  
415 plus the answering of the generation model. We 462  
416 employ LLaMA-3.1-Instruct with int4 quantization 463  
417 for compression. We choose HotpotQA, where 464  
418 the average context length is about tens of thou- 465  
419 sands tokens. As shown in Table 4, the overall time 466  
420 cost is significantly reduced compared to baseline 467  
421 methods, and the quantized model can also ensure 468  
422 similar performance on results. 469

### 5.3 Case Study 470

Table 1 presents a practical example of Attention- 471  
Prompt, demonstrating how our method effectively 472  
compresses the context while maintaining both ac- 473  
curacy and readability. Unlike the original con- 474  
text, which contains redundant information that 475  
can negatively affect the response quality, Atten- 476  
tionPrompt generates a concise and coherent output 477  
with minimal token usage. Other methods, such 478  
as LLMLingua2, while providing a more compact 479  
result, produce fragmented and less readable re- 480  
sponses that lose coherence and relevance. Simi- 481  
larly, LongLLMLingua, despite reducing the con- 482  
text, fails to provide a clear and focused answer. 483  
In contrast, AttentionPrompt generates the correct 484  
answer, “Ozalj,” with the highest compression ra- 485  
tio, illustrating its ability to preserve the essential 486  
information. This highlights AttentionPrompt’s 487  
capacity to enhance overall response quality, effec- 488  
tively balancing compression and clarity without 489  
introducing unnecessary complexity. 490

### 5.4 Ablation Study 491

**Fixed Hint Prefix** As discussed in Section Con- 492  
struct Answer Hint Prefix, we use a fixed hint prefix 493  
for queries that cannot generate one dynamically. 494  
The summarization performance under this setting 495  
is reported in Table 2. To further validate the robust- 496  
ness of this approach, we conduct experiments on 497  
three additional datasets: 2WikiMQA, HotpotQA, 498  
and TriviaQA. As shown in Table 5, we experiment 499  
with LLaMA-3.1-8B-Instruct, using only a fixed 500  
hint prefix results in a slight performance drop com- 501  
pared to our original method, yet still consistently 502  
outperforms all baselines. This demonstrates that 503  
even a static hint prefix serves as an effective an- 504  
chor token, reliably guiding the model to identify 505  
important information within the context. 506

**Combination with other RAG methods** To ex- 507  
plore the integration potential of our method with 508  
retrieval-based techniques, we implemented a two- 509  
stage pipeline: we first use vanilla RAG (BGE-M3) 510  
for initial context retrieval, followed by our method 511  
for compression. Specifically, we retain the RAG 512  
chunk size at 300 tokens, while retrieve the top 15 513  
relevant chunks, and then apply our compression 514  
method. We use LLaMA-3.1-8B-Instruct as the 515  
generation model and conduct experiments on the 516  
HotpotQA, 2WikiMQA, and TriviaQA benchmarks 517  
(See Table 6). 518

**Hyperparameters** We conduct ablation studies on 519

Model	Method	2WikiMQA			HotpotQA			TriviaQA			Gov	
		EM↑	LLM Judge↑	CR↑	EM↑	LLM Judge↑	CR↑	EM↑	LLM Judge↑	CR↑	Rouge↑	CR↑
Llama-3.1-8B-Instruct	Original Prompt	0.49	0.46	1.0x	0.50	0.62	1.0x	0.87	0.80	1.0x	0.32	1.0x
	Random	0.21	0.15	5.0x	0.05	0.01	3.3x	0.76	0.73	1.7x	–	–
	LLMLingua2 (small)	0.24	0.26	5.0x	0.40	0.45	3.3x	0.86	<b>0.82</b>	1.7x	0.23	3.3x
	LLMLingua2 (large)	0.33	0.30	5.0x	0.44	0.51	3.3x	0.84	0.76	1.7x	0.24	3.3x
	LongLLMLingua	0.23	0.22	3.3x	0.33	0.40	3.3x	0.83	0.72	1.7x	0.23	3.1x
	BGE-M3	0.40	0.35	3.5x	0.42	0.50	3.0x	0.80	0.77	1.9x	0.21	3.0x
	AttentionPrompt (small)	0.39	0.36	6.3x	0.45	0.55	3.7x	0.84	0.77	2.0x	<b>0.28</b>	3.4x
	AttentionPrompt (large)	<b>0.42</b>	<b>0.38</b>	15x	<b>0.48</b>	<b>0.61</b>	5.6x	<b>0.89</b>	<u>0.81</u>	1.7x	<b>0.28</b>	4.2x
Qwen2.5-7B-Instruct	Original Prompt	0.56	0.46	1.0x	0.58	0.68	1.0x	0.94	0.82	1.0x	0.31	1.0x
	Random	0.18	0.06	5.0x	0.04	0.01	3.3x	0.65	0.51	1.7x	–	–
	LLMLingua2 (small)	0.29	0.27	5.0x	0.43	0.47	3.3x	<b>0.91</b>	0.80	1.7x	0.21	3.3x
	LLMLingua2 (large)	0.40	0.28	5.0x	0.48	<b>0.58</b>	3.3x	0.87	<b>0.81</b>	1.7x	0.22	3.3x
	LongLLMLingua	0.31	0.18	4.0x	0.28	0.31	3.3x	0.83	0.73	1.7x	0.23	3.1x
	BGE-M3	0.40	0.35	3.5x	0.42	0.50	3.0x	0.80	0.77	1.9x	0.22	3.0x
	AttentionPrompt (small)	<b>0.41</b>	0.28	6.3x	0.51	0.54	3.7x	0.83	0.73	2.0x	<b>0.28</b>	3.4x
	AttentionPrompt (large)	<b>0.41</b>	<b>0.30</b>	15x	<b>0.53</b>	<u>0.56</u>	5.6x	<u>0.88</u>	<u>0.80</u>	1.7x	<b>0.28</b>	4.2x

Table 2: Performance comparison on QA datasets (2WikiMQA, HotpotQA, TriviaQA) and the Gov summarization dataset.

Model	Method	1K			2K			4K		
		EM↑	LLM Judge↑	CR↑	EM↑	LLM Judge↑	CR↑	EM↑	LLM Judge↑	CR↑
Llama-3.1-8B-Instruct	Original Prompt	0.74	0.64	1.0x	0.67	0.57	1.0x	0.71	0.60	1.0x
	Random	0.11	0.04	2.9x	0.09	0.05	3.3x	0.05	0.04	3.3x
	LLMLingua2 (small)	0.55	0.41	2.9x	0.38	0.29	3.3x	0.36	0.29	3.3x
	LLMLingua2 (large)	0.69	0.51	2.9x	0.53	0.41	3.3x	0.49	0.37	3.3x
	LongLLMLingua	0.63	0.59	2.9x	0.43	0.38	3.3x	0.47	0.36	3.3x
	AttentionPrompt (small)	<b>0.74</b>	<b>0.65</b>	3.8x	<b>0.58</b>	<b>0.49</b>	3.8x	<b>0.53</b>	<b>0.41</b>	5.6x
Qwen2.5-7B-Instruct	Original Prompt	0.87	0.78	1.0x	0.75	0.71	1.0x	0.80	0.75	1.0x
	Random	0.04	0.04	2.9x	0.05	0.03	3.3x	0.04	0.01	3.3x
	LLMLingua2 (small)	0.49	0.38	2.9x	0.37	0.30	3.3x	0.34	0.24	3.3x
	LLMLingua2 (large)	0.70	0.52	2.9x	0.55	0.41	3.3x	0.50	0.38	3.3x
	LongLLMLingua	0.54	0.47	2.9x	0.35	0.29	3.3x	0.37	0.29	3.3x
	AttentionPrompt (small)	<b>0.88</b>	<b>0.76</b>	3.8x	<b>0.62</b>	<b>0.57</b>	3.8x	<b>0.66</b>	<b>0.59</b>	5.6x

Table 3: Comparison of different methods on BABILong.

two hyperparameters: chunk size (Section Method) and the number of Top- $K$  tokens used to select sentences within each chunk (Section Compress with Attention). For long contexts (e.g., LongBench), we use larger chunks and higher  $K$  to handle dispersed information; for shorter contexts (e.g., Babilong), we use smaller chunks and lower  $K$  for finer granularity.<sup>2</sup> Using TriviaQA, we test various configurations and observe (Table 7) that increasing chunk size and reducing  $K$  lowers the compression ratio with minimal performance drop. This suggests that dynamic hyperparameter tuning offers a better trade-off between efficiency and accuracy than fixed-ratio compression.

**Size of Compression Models** To evaluate the model-size sensitivity of our compression approach, we compare the performance of AttentionPrompt when using LLaMA-3.1-Instruct 8B versus

<sup>2</sup>Details in Appendix.

70B to compute attention scores. As shown in Table 2, the performance gap is minimal, indicating that even a lightweight model can effectively capture the attention patterns needed for compression. This suggests that AttentionPrompt is largely agnostic to model size and can maintain strong performance without relying on large-scale models—offering substantial efficiency benefits. Furthermore, as demonstrated in Section Efficiency Analysis, our method can be further accelerated through quantization. Overall, AttentionPrompt is highly scalable and future-proof, with the potential to continuously benefit from advances in foundation model development while remaining efficient and adaptable.

**Attention Layer Choice** In LLMs, attention layers capture different levels of information—shallow layers focus on syntax, while deeper ones encode semantics (Ben-Artzy and Schwartz, 2024; Jin

Method	EM $\uparrow$	LLM Judge $\uparrow$	CR $\uparrow$	Latency	
				Compression	Answering
Original Prompt	0.50	0.62	1.0x	–	19.36
Vanilla RAG (BGE-M3)	0.42	0.50	3.0x	6.11	5.99
LLMLingua2 (small)	0.40	0.45	3.3x	0.44	5.39
LLMLingua2 (large)	0.44	0.51	3.3x	0.88	5.43
LongLLMLingua	0.42	0.52	3.3x	7.12	5.41
<b>AttentionPrompt</b>	<b>0.45</b>	<b>0.54</b>	<b>4.0x</b>	<b>3.99</b>	<b>4.94</b>

Table 4: Time latency for inference by various methods on HotpotQA

Dataset	EM $\uparrow$	LLM Judge $\uparrow$	CR $\uparrow$
HotpotQA	0.44	0.53	3.3x
2WikiMQA	0.44	0.37	5.2x
TriviaQA	0.81	0.75	2.0x

Table 5: Performance of AttentionPrompt using fixed hint prefix across datasets.

EM	LLM	Judge	CR
HotpotQA	0.46	0.58	9.1x (+5.4x)
2wikimqa	0.43	0.40	7.9x (+1.6x)
triviaqa	0.88	0.83	3.2x (+1.5x)

Table 6: Combination with RAG methods

et al., 2025). To fully exploit this hierarchy, we aggregate attention scores across all layers for compression, as detailed in Section [Compute Attention Features from LLM](#). This mitigates layer-specific bias and captures a broader information spectrum. We compare this approach with using shallow, middle, or deep layers alone on HotpotQA with Llama-3.1-8B-Instruct.

As shown in Table 8, full-layer aggregation yields superior performance, validating our strategy.

## 6 Related Work

To reduce the cost of long-context generation, both soft and hard prompt compression methods have been proposed. Soft prompt methods include Gist (Mu et al., 2024) and 500x Compressor (Li et al., 2024), which compress context into dense tokens with minimal loss. Hard prompt approaches like LLMLingua (Jiang et al., 2023), LongLLMLingua (Jiang et al., 2024), and LLMLingua2 (Pan et al., 2024) use token-level filtering, achieving substantial speedups and compression. In contrast, AttentionPrompt enhances LLM performance by selecting explainable, model-attended content with-

Chunk Size	Top-K	EM $\uparrow$	LLM Judge $\uparrow$	CR $\uparrow$
300	5	0.80	0.73	<b>3.2x</b>
300	10	0.84	0.75	2.2x
300	15	<b>0.87</b>	<b>0.77</b>	1.9x
100	10	<b>0.90</b>	<b>0.79</b>	1.9x
200	10	0.88	0.77	2.2x
400	10	0.85	0.77	<b>2.4x</b>

Table 7: Performance on TriviaQA with different chunk sizes and top-K values

Layer Subsets	EM $\uparrow$	LLM Judge $\uparrow$	CR $\uparrow$
0 - 10	0.35	0.43	<b>4.5x</b>
11 - 20	0.38	0.50	3.6x
21 - 31	0.40	0.48	3.7x
0 - 31	<b>0.42</b>	<b>0.54</b>	3.6x

Table 8: Performance on HotpotQA with different subset of layers

out retraining. We compare against LongLLMLingua and LLMLingua2 to demonstrate its efficiency and robustness.

## 7 Conclusion

In this paper, we propose AttentionPrompt, a novel attention-guided context pruning method for LLMs. The core part of our method is the formatted attention focus mechanism, which constructs an answer hint prefix and utilize a fixed hint prefix in a *next-token-prediction* format for each query, guiding the LLM to attend relevant tokens in the retrieved context through one token. We conduct extensive experiments on the 2WikiMQA, HotpotQA, TriviaQA, GovReport, and Babilong benchmarks, demonstrating its strong performance.

## 549 Limitation

550 In this section, we faithfully discuss the current lim-  
551 itations and potential avenues for future research.

552 Regarding the attention feature computation, we  
553 currently aggregate attention scores across all lay-  
554 ers. However, we believe this process can be opti-  
555 mized using more sophisticated algorithms to im-  
556 prove efficiency.

557 Additionally, while we propose a dynamic com-  
558 pression ratio, we have not yet developed methods  
559 for explicitly controlling or instructing the desired  
560 ratio. Determining and setting precise parameters  
561 to achieve a specific compression ratio is a chal-  
562 lenging task. In future work, we aim to investigate  
563 ways to provide more flexible and accurate control  
564 over the compression ratio.

## 565 References

566 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu,  
567 Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao  
568 Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang,  
569 and Juanzi Li. 2024. [Longbench: A bilingual, mul-  
570 titask benchmark for long context understanding.](#)  
571 *Preprint*, arXiv:2308.14508.

572 Amit Ben-Artzy and Roy Schwartz. 2024. [Attend first,  
573 consolidate later: On the importance of attention in  
574 different llm layers.](#) *Preprint*, arXiv:2409.03621.

575 Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu  
576 Lian, and Zheng Liu. 2024. [Bge m3-embedding:  
577 Multi-lingual, multi-functionality, multi-granularity  
578 text embeddings through self-knowledge distillation.](#)  
579 *Preprint*, arXiv:2402.03216.

580 Xin Cheng, Xun Wang, Xingxing Zhang, Tao Ge, Si-  
581 Qing Chen, Furu Wei, Huishuai Zhang, and Dongyan  
582 Zhao. 2024. [xrag: Extreme context compression  
583 for retrieval-augmented generation with one token.](#)  
584 *arXiv preprint arXiv:2405.13792*.

585 David Chiang and Peter Cholak. 2022. [Overcoming a  
586 theoretical limitation of self-attention.](#) In *Proceed-  
587 ings of the 60th Annual Meeting of the Association  
588 for Computational Linguistics (Volume 1: Long Pa-  
589 pers)*, pages 7654–7664, Dublin, Ireland. Association  
590 for Computational Linguistics.

591 Kevin Clark, Urvashi Khandelwal, Omer Levy, and  
592 Christopher D. Manning. 2019. [What does BERT  
593 look at? an analysis of BERT’s attention.](#) In *Pro-  
594 ceedings of the 2019 ACL Workshop BlackboxNLP:  
595 Analyzing and Interpreting Neural Networks for NLP*,  
596 pages 276–286, Florence, Italy. Association for Com-  
597 putational Linguistics.

598 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,  
599 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,  
600 Akhil Mathur, Alan Schelten, Amy Yang, Angela

Fan, and 1 others. 2024. [The llama 3 herd of models.](#) 601  
*Preprint*, arXiv:2407.21783. 602

OpenAI et al. 2024. [Gpt-4 technical report.](#) *Preprint*, 603  
arXiv:2303.08774. 604

Xiaodong Gu, Zhaowei Zhang, Sang-Woo Lee,  
Kang Min Yoo, and Jung-Woo Ha. 2022. [Contin-  
uous decomposition of granularity for neural para-  
phrase generation.](#) In *Proceedings of the 29th Inter-  
national Conference on Computational Linguistics*,  
pages 6369–6378, Gyeongju, Republic of Korea. In-  
ternational Committee on Computational Linguistics. 605  
606  
607  
608  
609  
610  
611

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara,  
and Akiko Aizawa. 2020. [Constructing a multi-hop  
qa dataset for comprehensive evaluation of reasoning  
steps.](#) In *Proceedings of the 28th International Con-  
ference on Computational Linguistics*, pages 6609–  
6625. 612  
613  
614  
615  
616  
617

Cheng-Yu Hsieh, Yung-Sung Chuang, Chun-Liang Li,  
Zifeng Wang, Long T. Le, Abhishek Kumar, James  
Glass, Alexander Ratner, Chen-Yu Lee, Ranjay Kr-  
ishna, and Tomas Pfister. 2024. [Found in the middle:  
Calibrating positional attention bias improves long  
context utilization.](#) *Preprint*, arXiv:2406.16008. 618  
619  
620  
621  
622  
623

Jie Huang and Kevin Chen-Chuan Chang. 2023. [To-  
wards reasoning in large language models: A survey.](#)  
*Preprint*, arXiv:2212.10403. 624  
625  
626

Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng  
Ji, and Lu Wang. 2021. [Efficient attentions for long  
document summarization.](#) In *Proceedings of the 2021  
Conference of the North American Chapter of the  
Association for Computational Linguistics: Human  
Language Technologies*, pages 1419–1436, Online.  
Association for Computational Linguistics. 627  
628  
629  
630  
631  
632  
633

Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing  
Yang, and Lili Qiu. 2023. [Llmlingua: Compressing  
prompts for accelerated inference of large language  
models.](#) *Preprint*, arXiv:2310.05736. 634  
635  
636  
637

Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng  
Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2024. [Longllmlingua: Accelerating and enhancing llms  
in long context scenarios via prompt compression.](#)  
*Preprint*, arXiv:2310.06839. 638  
639  
640  
641  
642

Mingyu Jin, Qinkai Yu, Jingyuan Huang, Qingcheng  
Zeng, Zhenting Wang, Wenyue Hua, Haiyan Zhao,  
Kai Mei, Yanda Meng, Kaize Ding, Fan Yang, Meng-  
nan Du, and Yongfeng Zhang. 2025. [Exploring  
concept depth: How large language models acquire  
knowledge and concept at different layers?](#) *Preprint*,  
arXiv:2404.07066. 643  
644  
645  
646  
647  
648  
649

Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke  
Zettlemoyer. 2017. [Triviaqa: A large scale distantly  
supervised challenge dataset for reading comprehen-  
sion.](#) In *Proceedings of the 55th Annual Meeting of  
the Association for Computational Linguistics (Vol-  
ume 1: Long Papers)*, pages 1601–1611. 650  
651  
652  
653  
654  
655

656	Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. 2024. <a href="#">Babilong: Testing the limits of llms with long context reasoning-in-a-haystack</a> . <i>Preprint</i> , arXiv:2406.10149.	710
657		711
658		712
659		713
660		714
661	Zongqian Li, Yixuan Su, and Nigel Collier. 2024. <a href="#">500xcompressor: Generalized prompt compression for large language models</a> . <i>Preprint</i> , arXiv:2408.03094.	715
662		716
663		
664		
665	Jesse Mu, Xiang Lisa Li, and Noah Goodman. 2024. <a href="#">Learning to compress prompts with gist tokens</a> . <i>Preprint</i> , arXiv:2304.08467.	
666		
667		
668	Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor Rühle, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and Dongmei Zhang. 2024. <a href="#">Llmlingua-2: Data distillation for efficient and faithful task-agnostic prompt compression</a> . <i>Preprint</i> , arXiv:2403.12968.	
669		
670		
671		
672		
673		
674	Haoran Que, Feiyu Duan, Liqun He, Yutao Mou, Wangchunshu Zhou, Jiaheng Liu, Wenge Rong, Zekun Moore Wang, Jian Yang, Ge Zhang, Junran Peng, Zhaoxiang Zhang, Songyang Zhang, and Kai Chen. 2024. <a href="#">Hellobench: Evaluating long text generation capabilities of large language models</a> . <i>Preprint</i> , arXiv:2409.16191.	
675		
676		
677		
678		
679		
680		
681	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <a href="#">Squad: 100,000+ questions for machine comprehension of text</a> . <i>Preprint</i> , arXiv:1606.05250.	
682		
683		
684		
685	Yunxiao Shi, Xing Zi, Zijing Shi, Haimin Zhang, Qiang Wu, and Min Xu. 2024. <a href="#">Enhancing retrieval and managing retrieval: A four-module synergy for improved quality and efficiency in rag systems</a> . <i>Preprint</i> , arXiv:2407.10670.	
686		
687		
688		
689		
690	Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. 2023. <a href="#">Max-margin token selection in attention mechanism</a> . <i>Preprint</i> , arXiv:2306.13596.	
691		
692		
693		
694	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. <a href="#">Attention is all you need</a> . <i>Preprint</i> , arXiv:1706.03762.	
695		
696		
697		
698	Sourav Verma. 2024. <a href="#">Contextual compression in retrieval-augmented generation for large language models: A survey</a> . <i>arXiv preprint arXiv:2409.13385</i> .	
699		
700		
701	Minzheng Wang, Longze Chen, Fu Cheng, Shengyi Liao, Xinghua Zhang, Bingli Wu, Haiyang Yu, Nan Xu, Lei Zhang, Run Luo, Yunshui Li, Min Yang, Fei Huang, and Yongbin Li. 2024. <a href="#">Leave no document behind: Benchmarking long-context LLMs with extended multi-doc QA</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 5627–5646, Miami, Florida, USA. Association for Computational Linguistics.	
702		
703		
704		
705		
706		
707		
708		
709		
	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jixi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 23 others. 2025. <a href="#">Qwen2.5 technical report</a> . <i>Preprint</i> , arXiv:2412.15115.	716
	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. <a href="#">Hotpotqa: A dataset for diverse, explainable multi-hop question answering</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2369–2380.	717
		718
		719
		720
		721
		722
		723
	Yun Zhu, Jia-Chen Gu, Caitlin Sikora, Ho Ko, Yinxiao Liu, Chu-Cheng Lin, Lei Shu, Liangchen Luo, Lei Meng, Bang Liu, and Jindong Chen. 2024. <a href="#">Accelerating inference of retrieval-augmented generation via sparse context selection</a> . <i>Preprint</i> , arXiv:2405.16178.	724
		725
		726
		727
		728
		729

## A Pseudocode

The Pseudocode of AttentionPrompt is provided in Algorithm 1. The algorithm takes a retrieved long context, query, and two language models as input. First, it generates an answer hint prefix based on the query to guide the attention mechanism. Then, it splits the long context into chunks of size  $m$ . For each chunk, it generates an anchor token using the compression model. If the anchor token is valid (not “none”), it computes attention features using the anchor token and compresses the chunk accordingly. Finally, all compressed chunks are concatenated and used with the original query to generate the final answer.

---

### Algorithm 1 Pseudocode of AttentionPrompt.

---

```
1: Input: long context  $C$ , query  $q$ , generation
   model  $L$ , compression model  $L_C$ 
2: Output: Generated response  $y$ 
3:
4: Generate Answer Hint Prefix
5: Get answer hint prefix  $p$  through  $L$  with  $q$ 
6: Chunking
7: Generate chunks  $c_1, \dots, c_n$  by partitioning  $C$ 
   with chunk size  $m$ , where  $n = \lceil |C|/m \rceil$ 
8: Initialize empty variable  $C'$ 
9: Compressing with Attention
10: for  $j = 1$  to  $n$  do
11:   Generate the anchor token  $a_1$  with  $L_C$ ,  $c_j$ ,
      $q$ , and  $p$ 
12:   if  $a_1$  is "none" then
13:     continue
14:   else
15:     Obtain Attention Features  $A_1$  with  $a_1$ 
     and  $c_j$  ▷ Eq. (2)
16:     Get compressed  $c'_j$  according to  $A_1$  and
      $c_j$  ▷ Eq. (3)
17:     Append  $c'_j$  to  $C'$ 
18:   end if
19: end for
20: Generate  $y$  from  $L$  with  $C'$  and  $q$ 
21: Return  $y$ 
```

---

## B Implementation Details

### B.1 Prompt for generating answer prefix hint

We use the following prompt for generating answer prefix hint according to each query.

You are a formatting assistant. Given a query, your task is to generate a corresponding answering format. The format should maintain the same structure as the query but transform it into an incomplete answer template. If it is impossible to generate a format, return “None”.

The format is like a complete answer, but truncated before the keyword, and the keyword is not included in the format.

For instance, if the query is “Where is Daniel?”, the format should be “Daniel is in the”, as the next word is the key word.

Note: For yes/no questions, such as “Is Tom here?”, return “None” because these questions are typically answered with “yes” or “no” and do not have a natural continuation that leads to a single keyword.

Examples:

1. Question: Where is Daniel?

Format: Daniel is in the

2. Question: What time is it?

Format: It is

3. Question: Who is responsible for this?

Format: The person responsible for this is

4. Question: Which film was released more recently, Dance With A Stranger or Miley Naa Miley Hum?

Format: The film released more recently was

5. Question: Is Tom here?

Format: None

In generation , you should only return the format, not any other text.

Now, here’s a new question:

Question: question

Format:

### B.2 Fixed Hint Prefix

Our fixed hint prefix is: “The most relevant keyword or phrase to the context is \_\_\_\_\_”

### B.3 Prompt for generating anchor token

We use the following prompt to generate the anchor token for computing attention features.

You will be given a long context beginning with 'Context:', a question beginning with 'Question:', and a hint beginning with 'Hint:'. Please answer the question.

Context: {chunk}

Hint: You should answer begin with {prefix\_hint}, if there is no useful information in the context for the question in the context and you really don't know the answer, just answer prefix\_hint none.

Question: {question}

Answer:

{prefix\_hint}

### B.4 Hyperparameter settings

To reduce the randomness, we use greedy decoding in open-source LLMs generation. For the chunk size and  $K$  in the attention-based compression process, we set them according to the context length in different benchmarks. In LongBench, where contexts are quite long, we use larger chunk size and  $K$ , in contrast, in BABILong, where we choose to experiment with mid-sized context, we use smaller chunk size and  $K$ . The detailed setting is shown in Table 9.

Dataset	Chunk_Size	$K$
HotpotQA	300	12
2WikiMQA	300	15
TriviaQA	150	8
BABILong 1k	50	8
BABILong 2k	100	10
BABILong 4k	200	12

Table 9: Hyperparameter settings of the experiment