

ATTENTIONRAG: Attention-Guided Context Pruning in Retrieval-Augmented Generation

Anonymous ACL submission

Abstract

While RAG demonstrates remarkable capabilities in LLM applications, its effectiveness is hindered by the ever-increasing length of retrieved contexts, which introduces information redundancy and substantial computational overhead. Existing context 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 ATTENTIONRAG, an attention-guided context pruning method for RAG systems. The core idea of ATTENTIONRAG lies in its attention focus mechanism, which reformulates RAG 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 retrieved contexts. Extensive experiments on LongBench and Babilong benchmarks demonstrate that ATTENTIONRAG achieves up to 6.3x context compression while preserving or even exceeding the performance of uncompressed contexts.

1 Introduction

Retrieval-Augmented Generation (RAG) (Lewis et al., 2021) demonstrated remarkable capabilities in Large Language Model (LLM) 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). Although LLMs have demonstrated superior performance, they often lack domain knowledge. By directly leveraging external documents, RAG provides a lightweight, cost-efficient solution to expand the LLMs’ knowledge base without retraining their parameters.

While effective, RAG-based systems encounter significant challenges in handling long contexts. As the number of retrieved documents grows, the context becomes excessively long, introducing large

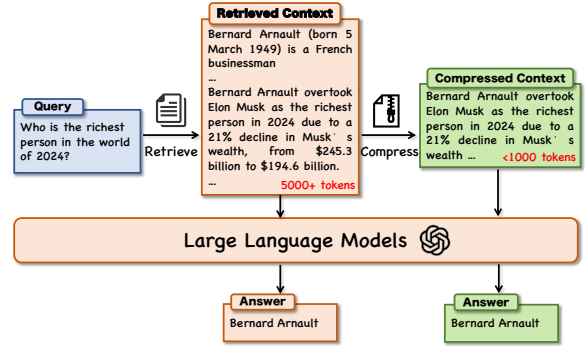


Figure 1: Illustration of RAG context compression

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). For example, LLM-Lingua (Jiang et al., 2023) compresses 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 during retrieval. This gap is primarily due to two fundamental challenges that hinder the effectiveness of attention in the retrieval architecture: (1) The long nature of RAG 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

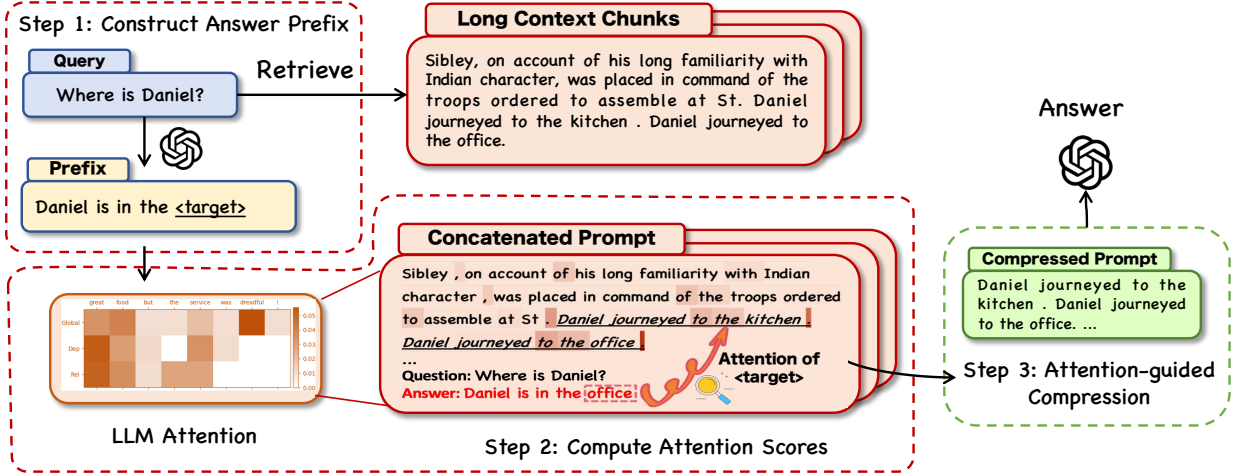


Figure 2: Illustration of ATTENTIONRAG

retrieved 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 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 ATTENTIONRAG (Attention-Guided Retrieval-Augmented Generation), which improves the relevance of the information extracted from the retrieved context by leveraging attention scores from intermediate layers of LLMs. To enable attention across queries and context, ATTENTIONRAG 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 retrieved 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. Next, we take the retrieved 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 ATTENTIONRAG 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 question-aware method for long-context pruning in RAG systems.
- We introduce a novel attention focus mechanism by reformulating RAG queries into a next-token prediction template, enabling precise and efficient computations of attention between queries and retrieved contexts.
- We conduct extensive experiments on LongBench and Babilong benchmarks. Results demonstrate the effectiveness of ATTENTIONRAG in long-context RAG systems.

2 Preliminaries

2.1 Retrieve-Augmented Generation

Retrieve-Augmented Generation (RAG) is a framework that enhances the capabilities of LLMs by integrating external knowledge through retrieval. A RAG system typically consists of two components: a *retriever*, which fetches relevant documents, called contexts, from a large corpus based

on a query, and a *generator*, which generates an answer using both the retrieved context and the model’s internal knowledge. This combination enables more accurate and contextually relevant outputs, especially for tasks requiring detailed or up-to-date information that might not be present in the model’s training data.

2.2 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 optimize RAG systems by compressing the retrieved context without compromising performance.

Problem Formulation Given a query q , a RAG system retrieves a relevant set of documents $\{d_1, \dots, d_K\}$ from a text corpus D . The retrieved documents are concatenated into a retrieved context $c = d_1 \oplus \dots \oplus d_N$. A large language model $P_\phi(a|q, c)$ takes the question and the retrieved context c as input and produces an answer a to the question. Our goal is to compress the retrieved context c into a dense one c' such that $|c'| \ll |c|$ while the LLM maintains the quality of generated answers when taking c' as input.

In this paper, we propose a novel attention-guided context pruning method called ATTENTION-RAG. The key idea of ATTENTIONRAG is re-

formulating each RAG 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 (§3.1); 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 (§3.2). Finally, we perform attention-guided compression: using attention scores from a language model, we identify and retain the most relevant parts of the retrieved context. Each of the steps is elaborated in the following sections.

3.1 Construct Answer Hint Prefix

To improve the alignment between the query and context, we associate each query with an answer hint prefix that allows the LLM to calculate the query-context attention through one 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 ____”. 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 related to the query.

We prompt the LLM to construct the answer hint prefix¹. 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, queries like wh-questions, the hint prefix can be derived by the query itself. In contrast, queries like yes/no-questions, where the answer’s first token is “Yes/No”, already align with

¹We use GPT-4o Mini for the generation, the prompt detail can be referred to §B.

the next-token-prediction paradigm. Leveraging the semantic understanding capabilities of LLMs, we prompt them with example answer hint prefix of various types questions 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 retrieved context C produced by the RAG framework 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, instruction (we instruct the LLM to generate "none" after hint prefix if the chunk is irrelevant, which will be used in §3.3) 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_i 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 §5.3.

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.

We provide the pseudocode of our method in Algorithm 1 in Appendix.

4 Experimental Setup

In this experiment, we evaluate the efficacy of ATTENTIONRAG 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. We evaluate ATTENTIONRAG in two key aspects: long context understanding and long context reasoning, reflected in two suites of benchmarks respectively:

LongBench (Bai et al., 2024): LongBench provides an extensive framework to evaluate an LLM’s ability to comprehend and process long-form texts. It encompasses a diverse suite of tasks—including complex information extraction, multi-step reasoning, and coherent text generation—that are designed to challenge models when key details are interspersed throughout extended passages. For our experiments, we specifically incorporate datasets such as TriviaQA (Joshi et al., 2017), HotpotQA (Yang et al., 2018), and 2WikiMQA (Ho et al., 2020), as each brings unique demands: 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. To ensure zero-shot consistency in testing, we do not use the example section from LongBench/TriviaQA.

BABILong (Kuratov et al., 2024): The BABILong benchmark provides a comprehensive framework for evaluating an LLM’s ability to reason over long contexts. It includes 20 varied tasks, such as fact chaining, simple induction, deduction, counting, and list/set manipulation, which become increasingly challenging when the relevant information is spread across extended natural texts. To assess ATTENTIONRAG’s reasoning capabilities, we select test splits of 1kqa1, 2kqa1, and 4kqa1 tokens to demonstrate its performance across different context lengths.

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 question, the model-generated answer, and the ground truth answer into GPT-4, asking it to deter-

mine whether the provided answer is correct. The prompt used for the scoring can be found in §B.

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 ATTENTIONRAG with the LLMLingua family, which includes LLMLingua (Jiang et al., 2023), LongLLMLingua (Jiang et al., 2024), and LLMLingua2 (Pan et al., 2024). Specifically, we choose LLMLingua2 and LongLLMLingua for baseline. They are state-of-the-art context compression methods. These baselines consist of question-aware methods (LongLLMLingua) and question-unaware methods (LLMLingua and LLMLingua2). The variety in baseline selection highlights the robustness of ATTENTIONRAG. To ensure a fair comparison, we maintain a compression rate for the LLMLingua family that is equal to or greater than that of ATTENTIONRAG. Additionally, we compare our results against the uncompressed context for reference, denoted as “Original Prompt”, and random compression, denoted as “Random”.

4.4 LLMs for Compression

We select two open-source LLMs for evaluation, Llama-3.1-8B-Instruct (Dubey et al., 2024) and Qwen-2.5-7B-Instruct (Yang et al., 2025). We experiment with both 8B and 70B versions of Llama-3.1-Instruct as our compression model, denoted as ATTENTIONRAG (8B) and ATTENTIONRAG (70B) respectively. Detailed hyperparameter configurations are provided in the Appendix §B.3.

5 Results

5.1 Overall Results

As the results in Tables 1 and 2 show, ATTENTIONRAG outperforms nearly all baseline methods across the benchmarks, consistently achieving superior results. On the LongBench benchmark, which includes contexts of tens of thousands of tokens, ATTENTIONRAG outperforms all LLMLingua methods across nearly all metrics. Specifically, in TriviaQA, ATTENTIONRAG outperforms the uncompressed method in EM score, demonstrating the effectiveness of our approach. However, the low compression ratio is due to two parallel factors: the scattering of relevant information across

Model	Method	2WikiMQA			HotpotQA			TriviaQA		
		EM↑	LLM Judge↑	CR↑	EM↑	LLM Judge↑	CR↑	EM↑	LLM Judge↑	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
	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	0.82	1.7x
	LLMLingua2 (large)	0.33	0.30	5.0x	0.44	0.51	3.3x	0.84	0.76	1.7x
	LongLLMLingua	0.23	0.22	3.3x	0.33	0.40	3.3x	0.83	0.72	1.7x
	ATTENTIONRAG (8B)	0.39	0.36	6.3x	0.45	0.55	3.7x	0.84	0.77	2.0x
	ATTENTIONRAG (70B)	0.42	0.38	15x	0.48	0.61	5.6x	0.89	<u>0.81</u>	1.7x
Qwen2.5-7B-Instruct	Original Prompt	0.56	0.46	1.0x	0.58	0.68	1.0x	0.94	0.82	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	0.91	0.80	1.7x
	LLMLingua2 (large)	0.40	0.28	5.0x	0.48	0.58	3.3x	0.87	0.81	1.7x
	LongLLMLingua	0.31	0.18	4.0x	0.28	0.31	3.3x	0.83	0.73	1.7x
	ATTENTIONRAG (8B)	0.41	0.28	6.3x	0.51	0.54	3.7x	0.83	0.73	2.0x
	ATTENTIONRAG (70B)	0.41	0.30	15x	0.53	<u>0.56</u>	5.6x	<u>0.88</u>	<u>0.80</u>	1.7x

Table 1: Performance of different methods in 2WikiMQA, HotpotQA and TriviaQA 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
	ATTENTIONRAG (8B)	0.74	0.65	3.8x	0.58	0.49	3.8x	0.53	0.41	5.6x
	ATTENTIONRAG (70B)	0.74	0.65	3.8x	0.58	0.49	3.8x	0.53	0.41	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
	ATTENTIONRAG (8B)	0.88	0.76	3.8x	0.62	0.57	3.8x	0.66	0.59	5.6x
	ATTENTIONRAG (70B)	0.88	0.76	3.8x	0.62	0.57	3.8x	0.66	0.59	5.6x

Table 2: Comparison of different methods on BABILong.

the context and the use of uniform hyperparameters. All these factors contribute to the lower ratio, which we aim to enhance. We discuss this in more detail in the Ablation Study.

Similar results can be observed in the BABILong benchmark, which involves a wider range of RAG context lengths. Our method consistently outperforms all LLMLingua methods while achieving the highest compression ratio, further demonstrating the effectiveness of ATTENTIONRAG in varying context sizes.

5.2 Case Study

Table 3 presents a practical example of ATTENTIONRAG, demonstrating how our method effectively compresses the context while maintaining both accuracy and readability. Unlike the original context, which contains redundant information that can negatively affect the response quality, AT-

TENTIONRAG generates a concise and coherent output with minimal token usage. Other methods, such as LLMLingua2, while providing a more compact result, produce fragmented and less readable responses that lose coherence and relevance. Similarly, LongLLMLingua, despite reducing the context, fails to provide a clear and focused answer. In contrast, ATTENTIONRAG generates the correct answer, ‘‘Ozalj,’’ with the highest compression ratio, illustrating its ability to preserve the essential information. This highlights ATTENTIONRAG’s capacity to enhance overall response quality, effectively balancing compression and clarity without introducing unnecessary complexity.

5.3 Ablation Study

Efficiency Analysis As discussed in §3, we divide the context into smaller chunks and use attention mechanisms for compression. To speed up the

Query:	Where was the wife of Francis I Rákóczi born?	Answer: Ozalj
Original Context:	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. ✗
Random:	drada., of. th-century French Baron Er of-R sister ofga bishopther arch of, not any socialoli,... (1400 tokens)	The text does mention it. ✗
LLMLingua2:	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. ✗
LongLLMLingua:	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. ✗
ATTENTIONRAG:	... Life Early years and family Ilona was born Ilona Zrínyi in Ozalj ... She was the daughter of Petar Zrinski, Ban (vicero) of Croatia, the niece of both Miklós Zrínyi and Fran Krsto Frankopan and the wife of Francis Rákóczi I ... (273 tokens)	Ozalj ✓

Table 3: Examples of compression results by various methods

process, we can design the prompt with an answer template that allows the overall generation time to be equal to the forward pass time of the compression model plus the compressed generation time of the generation model. Since a lightweight model is used for compression, the overall generation time is significantly reduced.

Attention Layer Choice In LLMs, attention mechanisms across different layers capture varying levels of information. Specifically, earlier layers often focus on syntactic structures, while deeper layers tend to capture more abstract semantic representations (Ben-Artzy and Schwartz, 2024). This hierarchical processing enables LLMs to effectively understand and generate complex language patterns (Jin et al., 2025). We adopt a comprehensive approach by aggregating attention scores across all layers for compression. As discussed in §3.2, this method ensures that the compression process leverages the full spectrum of information captured by the model and mitigate the potential attention bias across layers. We compare our aggregated approach using all layers with methods using different layer subsets (shallow, middle, and deep). To evaluate the effectiveness of these approaches, we conduct experiments on the HotpotQA dataset using the Llama-3.1-8B-Instruct model. As shown in Table 4, our method of summing attention scores across all layers outperforms the subset-layer approaches, supporting our hypothesis that leveraging the full spectrum of information captured across all layers provides better performance.

Hyperparameters We conduct ablation studies on two hyperparameters: the chunk size that we fur-

Layer Subsets	EM↑	LLM Judge ↑	CR↑
0 - 10	0.35	0.43	4.5x
11 - 20	0.38	0.50	3.6x
21 - 31	0.40	0.48	3.7x
0 - 31	0.42	0.54	3.6x

Table 4: Performance on HotpotQA with different subset of layers

ther partition the retrieved lengthy contexts (§3) and the number of Tok-K tokens for selecting the sentence in each chunk in the attention-based compression process (§3.3). For longer contexts (e.g., LongBench tasks), we increase the chunk size and use higher top-K values, as the information is more dispersed and larger chunks require a higher top-K to avoid internal bias. In contrast, for shorter contexts (e.g., Babilong tasks), we use smaller chunk sizes and lower top-K values to better capture the relevant information.² We conducted experiments on the TriviaQA dataset, separately testing various configurations of chunk sizes and K values. As shown in Table 5, we find that decreasing K and increasing the chunk size significantly reduce the compression ratio, while only slightly affecting the model’s performance. This is because these factors may decrease the relative context retrieved in each chunk. Therefore, in addition to a fixed compression ratio, we introduce more dynamic parameters to better control the balance between efficiency and performance.

²The hyperparameter settings for these experiments are detailed in §B.

Chunk Size	Top-K	EM↑	LLM Judge ↑	CR↑
300	5	0.80	0.73	3.2x
300	10	0.84	0.75	2.2x
300	15	0.87	0.77	1.9x
100	10	0.90	0.79	1.9x
200	10	0.88	0.77	2.2x
400	10	0.85	0.77	2.4x

Table 5: Performance on TriviaQA with different chunk sizes and top-K values

Size of Compression Models To assess the effectiveness of our lightweight compression model, we compare the results of compressing LLaMA-3.1-Instruct 8B against LLaMA-3.1-Instruct 70B. Our findings indicate that within the attention mechanism, the lightweight model delivers performance comparable to the larger model, suggesting that smaller models can achieve similar results. As shown in Table 1, the difference in performance between the large and small models is small, highlighting that when only attention scores are required, the smaller model is sufficiently capable for the compression task and can significantly enhance efficiency.

6 Related Work

Retrieval-Augmented Generation RAG has demonstrated effectiveness in various domains, such as open-domain question-answering (Han et al., 2024). However, the presence of redundant and irrelevant retrieved information can significantly hinder the performance of LLMs (Shi et al., 2023). To address this challenge, numerous methods have been proposed to improve the quality of retrieved information (Wang et al., 2023; Xu et al., 2023). For instance, Wang et al. (2023) presents a technique that trains models to identify and filter out irrelevant contexts in RAG systems, thereby enhancing both retrieval accuracy and downstream task performance. Xu et al. (2023), on the other hand, proposes a method that leverages context compression by training extraction models, which retain critical information while minimally impacting performance. In contrast to these approaches, ATTENTIONRAG relies on the inherent ability of LLMs to recognize useful information within their hidden states without requiring additional training, and can be applied to common compression problems. Our method not only achieves superior performance but is also more transferable and readily applicable across different models.

Prompt Compression To address the efficiency and cost challenges of long-context generation, various prompt compression methods have been proposed. These methods can be broadly categorized into soft prompts and hard prompts. In the domain of soft prompts, Mu et al. (2024) introduces the gist token, achieving a high compression rate with minimal performance loss. Similarly, Li et al. (2024) compresses long contexts into a single special token, retaining 50% to 70% of the original performance while significantly reducing input length. For hard prompts, Jiang et al. (2023) presents LLMingua, a method that identifies and removes unimportant tokens from prompts, optimizing inference efficiency. Building upon this, Jiang et al. (2024) proposes LongLLMingua, which further enhances the ability of language models to capture key information in long-context scenarios through prompt compression. This approach achieves up to a 17.1% performance improvement with a 4x compression ratio, effectively addressing the computational and latency challenges of processing long contexts. Additionally, Pan et al. (2024) introduces LLMingua-2, a task-agnostic prompt compression method trained via data distillation from GPT-4 for token classification with a BERT-level encoder. LLMingua-2 excels in handling out-of-domain data, offering 3x-6x faster performance compared to its predecessors. In contrast to these methods, ATTENTIONRAG focuses on improving LLM performance by leveraging explainable contexts derived from internal features. Our experiments will compare ATTENTIONRAG with LLMingua2 and LongLLMingua to highlight its effectiveness and efficiency.

7 Conclusion

In this paper, we propose ATTENTIONRAG, a novel attention-guided context pruning method for RAG systems. The core part of our method is the formatted attention focus mechanism, which constructs an answer hint prefix in a *next-token-prediction* format for each query, guiding the LLM to pay attention to relevant tokens in the retrieved context through one token. We conduct extensive experiments on the LongBench and Babilong benchmarks. Results show that ATTENTIONRAG achieves up to 6.3x context compression while preserving or even exceeding the performance of uncompressed contexts.

Limitation

In this section, we faithfully discuss the current limitations and potential avenues for future research.

Although we introduce the answer hint prefix paradigm, which can be effectively applied to a wide range of queries, there remain complex queries that it cannot fully address. For example, queries containing multiple sub-questions cannot be easily handled by locating a single focal token through the hint prefix. We propose that future work could explore query decomposition into smaller, unit queries, allowing for separate compression of each part.

Regarding the attention feature computation, we currently aggregate attention scores across all layers. However, we believe this process can be optimized using more sophisticated algorithms to improve efficiency.

Additionally, while we propose a dynamic compression ratio, we have not yet developed methods for explicitly controlling or instructing the desired ratio. Determining and setting precise parameters to achieve a specific compression ratio is a challenging task. In future work, we aim to investigate ways to provide more flexible and accurate control over the compression ratio.

References

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. 2024. [Longbench: A bilingual, multitask benchmark for long context understanding](#). *Preprint*, arXiv:2308.14508.

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

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

David Chiang and Peter Cholak. 2022. [Overcoming a theoretical limitation of self-attention](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7654–7664, Dublin, Ireland. Association for Computational Linguistics.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? an analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP*:

Analyzing and Interpreting Neural Networks for NLP, pages 276–286, Florence, Italy. Association for Computational Linguistics.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.

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

Xiaodong Gu, Zhaowei Zhang, Sang-Woo Lee, Kang Min Yoo, and Jung-Woo Ha. 2022. [Continuous decomposition of granularity for neural paraphrase generation](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6369–6378, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Rujun Han, Yuhao Zhang, Peng Qi, Yumo Xu, Jenyuan Wang, Lan Liu, William Yang Wang, Bonan Min, and Vittorio Castelli. 2024. [Rag-qa arena: Evaluating domain robustness for long-form retrieval augmented question answering](#). *Preprint*, arXiv:2407.13998.

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 Conference on Computational Linguistics*, pages 6609–6625.

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

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

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

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.

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

A Pseudocode

The Pseudocode of ATTENTIONRAG 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 ATTENTIONRAG.

```

1: Input: Retrieved long context  $C$ , query  $q$ , generation model  $L$ , compression model  $L_C$ 
2: Output: Generated sequence  $y$ 
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 Generated sequence  $y$ 

```

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 question, your task is to generate a corresponding answering format. The format should maintain the same structure as the question but transform it into an incomplete answer template. If it is impossible to generate a format, return "None".

The format is like an complete answer, but truncated before the key word, and the key word is not included in the format.

For instance, if the question 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 Prompt for generating anchor token

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

You will be given a long context begin with 'Context:', a question begin with 'Question:', and a hint begin 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.3 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 6.

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 6: Hyperparameter settings of the experiment