

# EXIT: Context-Aware Extractive Compression for Enhancing Retrieval-Augmented Generation

Anonymous ACL submission

## Abstract

We introduce **EXIT**, an extractive context compression framework that enhances both the effectiveness and efficiency of retrieval-augmented generation (RAG) in question answering (QA). Current RAG systems often struggle when retrieval models fail to rank the most relevant documents, leading to the inclusion of more context at the expense of latency and accuracy. While abstractive compression methods can drastically reduce token counts, their token-by-token generation process significantly increases end-to-end latency. Conversely, existing extractive methods reduce the latency but rely on independent, non-adaptive sentence selection, failing to fully utilize contextual information. EXIT addresses these limitations by classifying sentences from retrieved documents—while preserving their contextual dependencies—enabling parallelizable, context-aware extraction that adapts to query complexity and retrieval quality. Our evaluations on both single-hop and multi-hop QA tasks show that EXIT consistently surpasses existing compression methods and even uncompressed baselines in QA accuracy, while also delivering substantial reductions in inference time and token count. By improving both effectiveness and efficiency, EXIT provides a promising direction for developing scalable, high-quality QA solutions in RAG pipelines.<sup>1</sup>

## 1 Introduction

Retrieval-Augmented Generation (RAG) (Lewis et al., 2020b; Khandelwal et al., 2020) is the task of enhancing the responses of Large Language Models (LLMs) with relevant external contexts or documents. By grounding answers in evidence, RAG systems have gained much attention for mitigating hallucination issues (Ram et al., 2023; Li et al., 2023c) and improving factual reliability (Jeong et al., 2024; Xia et al., 2024b).

<sup>1</sup>We will make our code publicly available upon acceptance of this paper.

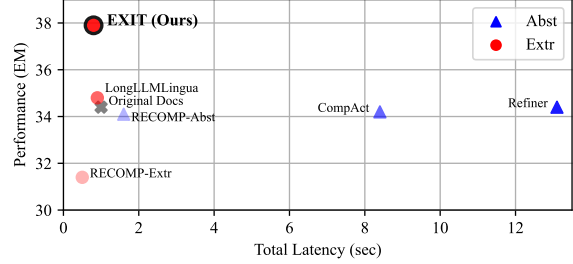


Figure 1: Average QA accuracy (EM) and efficiency (Total Latency) for various compression methods using Contriever-MSMARCO as the retriever and Llama-3.1-8b-Instruct as the reader.

However, they face significant challenges in practical deployment, as retrieval models sometimes fail to rank the most relevant documents at the top (Robertson and Zaragoza, 2009; Izacard et al., 2022). One potential solution is to retrieve a larger set of documents to ensure the coverage of the necessary information, but this approach compromises both effectiveness and efficiency. Specifically, for effectiveness, LLMs often struggle with processing long contexts, overlooking critical information located in the middle of contexts (Liu et al., 2024). Additionally, irrelevant information in retrieved documents can act as distractors, significantly degrading the overall QA performance (Shi et al., 2023; Li et al., 2023b; Wu et al., 2024). From an efficiency perspective, increasing the context size raises inference latency—due to quadratic complexity in attention computation (Xia et al., 2024a)—and API costs tied to input length<sup>2</sup>. Moreover, the context window limitations inherent in LLM architectures set strict upper bounds on the maximum input size (Peng et al., 2024).

To address these challenges, context compression has emerged as a promising solution, condensing essential information from multiple retrieved contexts through either abstractive or extractive approaches. While they can reduce inference time and filter out irrelevant information, both approaches

<sup>2</sup><https://openai.com/api/pricing/>

still have notable drawbacks. Specifically, abstractive compression methods—often implemented via autoregressive generation—summarize or rewrite documents into a single condensed passage (Li et al., 2023d; Xu et al., 2024; Yoon et al., 2024; Li et al., 2024; Wang et al., 2023), significantly increasing end-to-end latency due to their token-by-token generation process. For instance, as illustrated in Figure 1, CompAct (Yoon et al., 2024), a representative abstractive approach, takes over 8 seconds to process just five documents for a single query, whereas using the original document without compression takes only 1 second.

On the other hand, extractive compression approaches can offer a more efficient alternative (Xu et al., 2024; Jiang et al., 2023, 2024), by selecting relevant textual segments (e.g., sentences, or even token-level excerpts) directly from retrieved documents. This strategy reduces both compression time and overall latency. However, current extractive methods have yet to reach their full potential in terms of effectiveness (Choi et al., 2021; Pan et al., 2024). They often rely on rigid selection criteria that do not adapt well to variations in query complexity or the quality of retrieved documents, and they frequently neglect to fully leverage the broader context when choosing which tokens or sentences to retain. Specifically, as illustrated in Figure 1, while extractive approaches such as RECOMP-Extr (Xu et al., 2024) achieve minimal compression time, their inability to dynamically adjust selection processes results in suboptimal QA performance.

Therefore, in this work, we propose a novel compression framework for RAG, **EX**tractIve Context compression (EXIT), designed to enhance both effectiveness and efficiency by improving efficiency through an extractive compression and by enhancing effectiveness through dynamic, context-aware sentence selection. Specifically, as shown in Figure 2, EXIT operates in three stages: (1) splitting retrieved documents into sentences, (2) performing parallelizable binary classification (“Yes” or “No”) on each sentence to assess its relevance while considering its full document context, rather than evaluating them independently, and (3) recombining selected sentences while preserving their original order. Therefore, as shown in Figure 1, EXIT frames context compression as a sentence classification problem, enabling it to outperform both compression methods and the uncompressed baseline in terms of speed. Specifically, it reduces process-

ing time from several seconds to about 1 second. Moreover, by leveraging context-aware and adaptive sentence selection, EXIT also surpasses other extractive methods in accuracy. Also, we note that EXIT operates as a plug-and-play module that can be seamlessly integrated into any existing RAG pipeline without architectural modifications.

We evaluate EXIT on both single-hop QA tasks (NQ (Kwiatkowski et al., 2019), TQA (Joshi et al., 2017)) and multi-hop QA tasks (HQA (Yang et al., 2018), 2WikiMultiHopQA (Ho et al., 2020)). Experimental results demonstrate that EXIT not only improves effectiveness over both abstractive and extractive compression baselines but also significantly reduces latency compared to abstractive methods and the uncompressed baseline.

Our contributions are as follows:

- We identify and address the key weaknesses of existing context compression methods: abstractive methods incur prohibitive latency, while traditional extractive methods rely on rigid, non-adaptive content selection.
- We propose **EXIT** (**EX**tractIve Context compression), an extractive compression framework that dynamically adjusts to query complexity and retrieval quality.
- We demonstrate, through extensive experiments, that EXIT surpasses previous compression methods and uncompressed retrievals, improving QA performance while significantly reducing both token counts and end-to-end latency.

## 2 Related Work

**Retrieval-Augmented Generation.** RAG enhances LLMs by grounding responses in retrieved external documents (Lewis et al., 2020b; Shi et al., 2024; Ram et al., 2023). However, increasing retrieved content often leads to longer contexts, higher inference costs (Xia et al., 2024a), and retrieval noise (Shi et al., 2023; Liu et al., 2024), impacting efficiency and accuracy.

To mitigate these challenges, ad-hoc methods have been proposed at different stages of the retrieval pipeline, including pre-retrieval and post-retrieval methods. Pre-retrieval methods, such as Landmark Embedding (Luo et al., 2024) and Late Chunking (Günther et al., 2024), improve document representations to enhance retrieval effectiveness but lack adaptive filtering mechanisms after

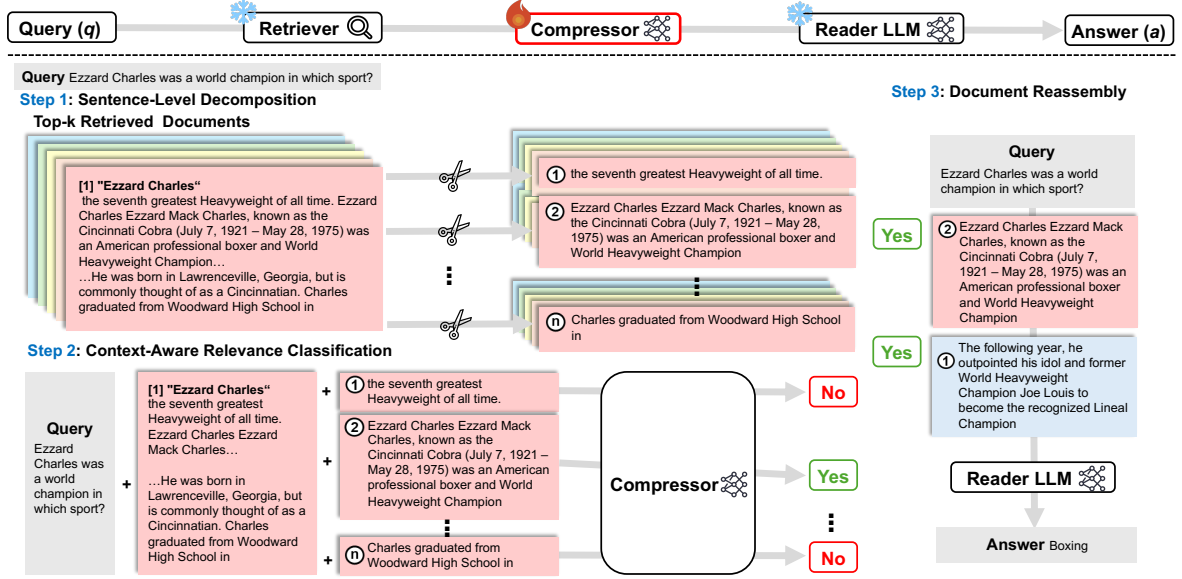


Figure 2: Overview of our framework. First, the retrieved document is split into sentences. Next, each sentence is classified as either “Yes” or “No” using the Compressor. Finally, sentences with scores above the threshold are recombined in their original order to complete the compression.

retrieval and may be difficult to integrate orthogonally with existing retrieval methods. Post-retrieval methods, including reranking (Nogueira and Cho, 2019; Qin et al., 2023; Li et al., 2023a) and context compression (Xu et al., 2024; Yoon et al., 2024; Li et al., 2024; Jiang et al., 2024), refine retrieved documents through reordering or pruning. However, most post-retrieval methods overlook query complexity and operate on a fixed number of retrieved items, limiting their adaptability in balancing effectiveness and efficiency.

We introduce EXIT, a novel post-retrieval compression framework that adaptively filters query-relevant sentences while discarding redundancy. Unlike existing methods, EXIT dynamically adjusts to query complexity and retrieval quality, balancing efficiency and effectiveness. Its lightweight, adaptive design seamlessly integrates into RAG pipelines while remaining orthogonal to pre-retrieval methods.

**Context Compression.** Context compression is a practical remedy for handling increasingly extensive prompt lengths in RAG pipelines. Existing methods commonly fall into soft or hard compression. Soft compression focuses on embedding prompts into compact continuous representations (Wingate et al., 2022; Mu et al., 2023; Ge et al., 2024; Chevalier et al., 2023; Cheng et al., 2024), but requires extensive training and architectural changes, making it unsuitable for black-box LLMs. Hard compression, by contrast, removes non-essential textual content directly (Li et al.,

2023e; Jiang et al., 2023), offering a plug-and-play solution even compatible with API-based models such as ChatGPT (OpenAI, 2023).

Hard compression are further divided into abstractive and extractive methods. Abstractive methods employ autoregressive models to generate query-focused summaries, thus reducing token counts at the cost of additional latency and potential hallucinations (Zhao et al., 2020). For example, RECOMP-Abst (Xu et al., 2024) uses a T5-based summarizer for token reduction but requires dataset-specific training and slows inference. CompAct (Yoon et al., 2024) and Refiner (Li et al., 2024) take this method further by leveraging even larger LLMs with 7B parameters, compounding latency issues and increasing resource demands. Extractive methods, such as RECOMP-Extr (Xu et al., 2024), directly select salient segments (e.g., sentences or tokens) from retrieved documents. This avoids the autoregressive bottleneck and mitigates hallucinations. However, their static and context-agnostic selection of only a few sentences per document limits its performance. Similarly, token-level methods such as LLMingua family (Li et al., 2023e; Jiang et al., 2023; Pan et al., 2024; Jiang et al., 2024) can distort semantic coherence by removing key entities or splitting essential facts.

In short, abstractive methods offer strong compression but suffer from latency and potential hallucinations, while extractive methods are often rigid and lack context awareness. Our work addresses these limitations by proposing a paralleliz-

able, context-aware extractive compression framework. It adaptively selects sentences at scale, optimizing accuracy and speed in complex retrieval scenarios.

### 3 Method

In this section, we first present our problem formulation, the RAG pipeline with a compression stage, and our novel compression framework, **EXIT**, which is designed to extract key evidence for answering in a parallel manner.

#### 3.1 Problem Formulation

**RAG Pipeline with Compression.** Given a query  $q$  and a document corpus  $\mathcal{C}$ , a RAG pipeline first retrieves Top- $k$  relevant document set  $D$ :

$$D = \{d_1, \dots, d_k\} = \text{Retriever}(q, \mathcal{C}), \quad (1)$$

The retrieved documents within the document set  $D$  are then processed by a compression module that preserves query-relevant information while significantly reducing input length:

$$D' = \text{Compressor}(q, D) \text{ s.t. } l(D') \ll l(D), \quad (2)$$

where  $l(\cdot)$  represents the function calculating the number of tokens in the document set. After compression, the number of tokens included in  $D'$  is substantially decreased compared to  $D$ . Finally, an LLM generates the answer  $a$  using the compressed set  $D'$  and the given query  $q$ :

$$a = \text{LLM}(q, D'). \quad (3)$$

**Objectives of Compression.** Effective compression in the RAG pipeline must satisfy three key criteria: (1) **Token Reduction**— $D'$  should have fewer tokens to speed up answer generation; (2) **Retention of Key Evidence**—essential information must be preserved to maintain answer accuracy; and (3) **Efficient Processing**—compression should be fast enough to avoid adding significant latency.

While token count reflects reading time, it alone cannot measure overall efficiency, as it ignores the computational cost of compression. A truly efficient pipeline minimizes both token count and processing overhead.

#### 3.2 Extractive Context Compression (EXIT)

To achieve three objectives of the compression step, EXIT consists of three main components: sentence-level decomposition, context-aware relevance classification, and document reassembly.

**Sentence-Level Decomposition.** In Step 1 of Figure 2, EXIT divides each retrieved document into individual sentences using a rule-based sentence tokenizer. For each document  $d_i \in D$ , it produces a sentence set  $S_i = \{s_{i1}, s_{i2}, \dots, s_{in}\}$ , where  $s_{ij}$  is the  $j$ -th sentence in document  $i$ . Operating at the sentence level avoids the fragmentation of key phrases and preserves entity relationships that token-level compression techniques (Jiang et al., 2023) often disrupt. As a result, the compressed context preserves both syntactic coherence and semantic integrity, ensuring that key information is retained.

**Context-Aware Relevance Classification.** To efficiently filter sentences in  $D$  that contain key evidence for answering a query, we introduce a relevance evaluation method based on **context awareness** and **single-token prediction**. First, incorporating the entire document  $d_i$  as context is essential, as understanding a sentence often requires the broader document context rather than an isolated sentence. This ensures that no relevant information is overlooked and enables more effective compression. Second, to maintain efficiency, we adopt lightweight single-token prediction instead of multi-token generation (Yoon et al., 2024; Li et al., 2024), which can introduce computational overhead. Inspired by Zhong et al. (2022), we leverage probabilistic classification for sentence relevance assessment. Given query  $q$ , document  $d_i$ , and sentence  $s_{ij}$ , the evaluation model predicts relevance using a binary classification with “Yes” and “No” labels:

$$r_{ij} = \frac{P(\text{“Yes”}|q, d_i, s_{ij})}{P(\text{“Yes”}|q, d_i, s_{ij}) + P(\text{“No”}|q, d_i, s_{ij})} \quad (4)$$

where  $P(\cdot|\cdot)$  represents the likelihood from the evaluation model. This computation is parallelized across sentences for efficiency.

EXIT then selects sentences with a relevance score above a predefined threshold  $\tau$ , ensuring that only key information is retained. Unlike fixed-size selection, our framework produces an **adaptive number of sentences** in the compressed set  $D'$ , aligning with prior work (Jeong et al., 2024) that acknowledges query-dependent information needs. This approach optimizes compression while preserving essential evidence.

**Document Reassembly.** As shown in Step 3 of Figure 2, EXIT reconstructs the compressed document  $D'$  by concatenating selected sentences in their original order, preserving coherence and logical



flow (Hwang et al., 2024) for accurate downstream reasoning.

### 3.3 Classifier Model Training

**Training Strategy.** Our goal is to train a relevance classifier capable of accurately identifying which sentences provide the evidence required to answer a query. To approximate real-world complexity, we utilize a question-answering dataset that requires multi-sentence reasoning and offers explicit sentence-level annotations of essential information. Leveraging these annotations, we model three typical retrieval outcomes: (1) sentences containing necessary evidence, (2) seemingly relevant sentences missing key details, and (3) entirely irrelevant sentences.

**Data Sampling.** To train a robust relevance classifier, we construct a diverse dataset reflecting various post-retrieval conditions. Positive samples include sentences essential for correct answers, while hard negatives consist of remaining sentences from the same documents. Additionally, random negatives pair queries with unrelated sentences, helping the model distinguish relevance from noise. This balanced sampling enables effective filtering of non-essential information without relying on explicit supervision.

**Training Procedure.** Each training instance is represented as  $(q, s, d, l)$ , where  $q$  is the query,  $s$  is a candidate sentence,  $d$  is the document containing  $s$ , and  $l \in \{\text{“Yes”}, \text{“No”}\}$  indicates whether  $s$  provides the required evidence. We employ a binary cross-entropy loss function to train the classifier:

$$\mathcal{L} = -\mathbb{1}_{l=\text{“Yes”}} \log P(\text{“Yes”}) - (1 - \mathbb{1}_{l=\text{“Yes”}}) \log P(\text{“No”}), \quad (5)$$

By exposing the classifier to a balanced and diverse set of retrieval scenarios, we improve its ability to generalize and reliably identify sentences that contain the critical evidence for answering queries.

## 4 Experiment Setups

We conduct comprehensive experiments to evaluate EXIT’s effectiveness and efficiency in context compression for RAG systems. More implementation details of our experiments are in Appendix B.

**Datasets.** We evaluate on both single-hop and multi-hop question answering datasets: **NaturalQuestions (NQ)** (Kwiatkowski et al., 2019) and **TriviaQA (TQA)** (Joshi et al., 2017) for single-hop QA; **HotpotQA (HQA)** (Yang et al., 2018) and **2WikiMultihopQA (2WIKI)** (Ho et al., 2020)

for multi-hop QA. We use the test set for TQA evaluations and development sets for all other datasets. For the train dataset for the classifier, we exploit the train split of HQA, which has relevant annotations to each sentence in the multiple documents for a query.

**Model Configuration.** Our system consists of three primary components. The retriever employs **Contriever-MSMARCO** (Izacard et al., 2022), a dense retriever fine-tuned on MSMARCO (Nguyen et al., 2016). The EXIT compressor utilizes **Gemma-2B-it** (Mesnard et al., 2024), optimized for efficient parallel processing. For the reader model, we exploit two scales of instruction-tuned models: **Llama3.1-{8, 70}B-Instruct** (Dubey et al., 2024).

**Baselines.** We compare EXIT against the following context compression methods: **1) Original Documents** serves as uncompressed baseline. For abstractive methods, **2) RECOMP-Abs** (Xu et al., 2024) uses T5-based (775M) summarization tuned for NQ, TQA, HQA (HQA model for 2WIKI), **3) CompAct** (Yoon et al., 2024) implements Mistral-7B-based iterative compression with 5-segment blocks, and **4) Refiner** (Li et al., 2024) uses Llama2-7B-based compression. For extractive methods, **5) RECOMP-Extr** (Xu et al., 2024) employs Contriever-based (110M) sentence-level extraction tuned for NQ, TQA, HQA (HQA model for 2WIKI), and **6) LongLLMLingua** (Jiang et al., 2024) uses Llama2-7B-chat for token-level extraction with 0.4 dynamic compression rate.

**Evaluation Metrics.** We use three metrics: **Exact Match (EM)** and **F1** score to measure effectiveness in question answering, and **end-to-end inference latency (Lat.)** in seconds to assess efficiency. Here, end-to-end latency measures the total time for compression and generation, providing a holistic efficiency metric beyond token count, which only reflects reading time.

**Implementation Details.** We conduct retrieval over the December 2018 Wikipedia dump and apply SpaCy for sentence splitting. We employ vLLM v0.5.5 (Kwon et al., 2023) for accelerated inference with the hyperparameter  $T = 0.0$  and  $\text{Top-}P = 1.0$ . We empirically set the relevance threshold  $\tau = 0.5$ . All experiments are conducted on A100-SXM4-80GB GPUs.

Table 1: Performance across models and datasets, measured by EM, F1, and inference latency (Lat.). 8B reader experiments were conducted on a single A100-80GB GPU, while 70B reader experiments utilized 4 A100-80GB GPUs in parallel. Best results for each dataset are highlighted in **bold**, and second best results are underlined. The “Type” column denotes whether a given compressor is abstractive (Abs.) or extractive (Ext.).

Compressor	Type	NQ			TQA			HQA			2WIKI			AVG.		
		EM ↑	F1 ↑	Lat. ↓	EM ↑	F1 ↑	Lat. ↓	EM ↑	F1 ↑	Lat. ↓	EM ↑	F1 ↑	Lat. ↓	EM ↑	F1 ↑	Lat. ↓
Llama3.1-8B-Instruct																
Original Docs	-	34.6	47.1	1.0	58.8	68.6	0.9	28.1	38.6	1.0	16.1	24.9	1.1	34.4	44.8	1.0
RECOMP-Abst	Abs.	31.3	43.2	1.6	55.9	65.7	1.4	26.5	37.0	2.2	<u>22.7</u>	<u>29.1</u>	2.1	34.1	43.7	1.8
CompAct	Abs.	32.9	44.6	8.5	58.1	67.7	8.8	<u>28.8</u>	39.8	8.3	16.8	26.0	8.1	34.2	44.5	8.4
Refiner	Abs.	32.9	45.0	28.1	59.2	<u>68.9</u>	10.9	<u>28.8</u>	<u>40.0</u>	6.9	16.8	25.4	6.4	34.4	<u>44.8</u>	13.1
RECOMP-Extr	Ext.	34.6	44.6	<b>0.5</b>	56.5	65.1	<b>0.4</b>	23.4	32.8	<b>0.4</b>	11.2	19.6	<b>0.6</b>	31.4	40.5	<b>0.5</b>
LongLLMLingua	Ext.	30.2	41.5	0.9	<u>59.4</u>	68.0	0.8	28.0	38.0	<u>0.8</u>	21.5	27.4	<u>0.9</u>	<u>34.8</u>	43.7	0.9
EXIT (Ours)	Ext.	<b>35.9</b>	<b>47.8</b>	<u>0.8</u>	<b>60.8</b>	<b>69.9</b>	<u>0.7</u>	<b>30.6</b>	<b>41.5</b>	<u>0.8</u>	<b>24.2</b>	<b>30.8</b>	<u>0.9</u>	<b>37.9</b>	<b>47.5</b>	<u>0.8</u>
Llama-3.1-70B-Instruct																
Original Docs	-	35.6	48.0	8.6	65.1	73.9	7.7	33.7	44.5	8.3	20.8	28.3	9.1	38.8	48.7	8.4
RECOMP-Abst	Abs.	34.1	47.0	4.5	61.3	70.6	3.3	30.3	40.8	4.4	24.2	30.3	4.2	37.5	47.2	4.1
CompAct	Abs.	34.1	45.4	11.9	62.6	71.1	11.7	33.8	44.1	11.0	20.5	27.4	11.6	37.8	47.0	11.5
Refiner	Abs.	35.3	<u>47.1</u>	42.5	64.3	73.0	18.3	33.8	44.7	14.6	21.2	28.0	11.2	38.7	48.2	21.6
RECOMP-Extr	Ext.	<u>35.8</u>	45.3	<b>2.5</b>	63.5	71.0	<b>2.2</b>	27.6	36.7	<b>2.9</b>	13.8	19.3	<b>3.3</b>	35.2	43.1	<b>2.7</b>
LongLLMLingua	Ext.	32.2	44.0	4.4	<u>66.7</u>	<u>75.2</u>	3.9	<u>34.1</u>	<u>45.3</u>	4.0	<u>28.3</u>	<b>34.8</b>	4.3	<u>40.3</u>	<u>49.8</u>	4.1
EXIT (Ours)	Ext.	<b>36.9</b>	<b>49.4</b>	<u>3.9</u>	<b>67.3</b>	<b>75.9</b>	<u>3.1</u>	<b>37.0</b>	<b>48.3</b>	<u>3.3</u>	<b>28.6</b>	<u>34.5</u>	<u>3.5</u>	<b>42.5</b>	<b>52.0</b>	<u>3.5</u>

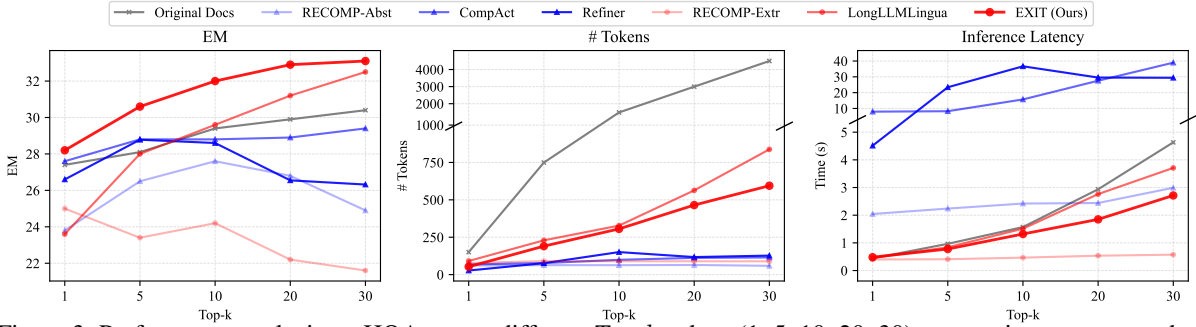


Figure 3: Performance analysis on HQA across different Top- $k$  values (1, 5, 10, 20, 30), comparing accuracy, token retention, and inference latency between baselines and our method. All experiments were conducted on a single A100-80GB GPU.

## 5 Main Results

Table 1 summarizes our evaluation results across multiple datasets and compression methods. With the 8B reader, EXIT demonstrates strong generalization: although trained solely on HQA, it effectively addresses both single-hop (NQ, TQA) and multi-hop (2WIKI) queries under out-of-domain conditions. Compared to all baseline methods, EXIT consistently improves EM scores—for instance, by 1.3 and 2.0 points on NQ and TQA, and by even larger margins of 2.5 and 8.1 points on HQA and 2WIKI, respectively. Notably, these accuracy gains come with an average latency of just 0.8s, substantially faster than abstractive compression methods.

The benefits of EXIT become more pronounced at larger scales. Using the 70B reader, EXIT surpasses the accuracy of all competing methods, averaging a 3.7-point improvement in EM and a 3.3-point improvement in F1 over the uncom-

pressed baseline. On HQA, it achieves a 3.3-point EM gain while maintaining an efficient 3.5s latency—faster than using uncompressed documents and still competitive with the previously fastest method, RECOMP-Extr, but with significantly higher accuracy. EXIT’s effectiveness and efficiency, especially with larger models, makes it a practical solution for large-scale QA applications.

## 6 Analyses

We conduct a series of analyses examining EXIT’s robustness, classification performance, latency factors, and design choices under various configurations. Additional experimental results and analyses are provided in Appendix C.

### 6.1 Robustness Analysis

To examine EXIT’s robustness as the retrieval set size grows, we gradually increased the number of retrieved documents ( $k \in \{1, 5, 10, 20, 30\}$ ) with an 8B reader, as shown in Figure 3. We found that

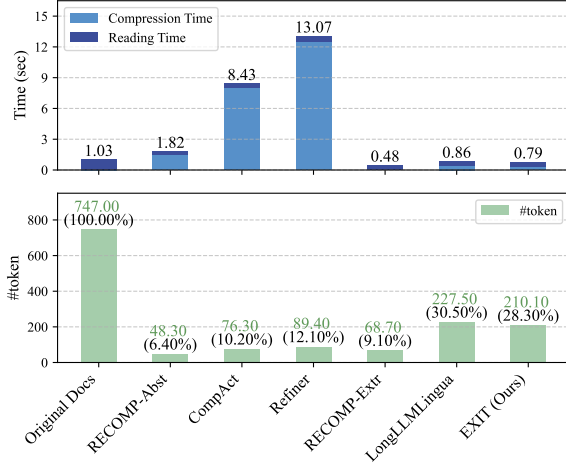


Figure 4: Comparison of compression and reading latency across baselines and our method in QA setting. Experiments were conducted on a single A100 GPU.

EXIT steadily improves EM scores—from 28.2 points at  $k = 1$  to 33.1 points at  $k = 30$ —while avoiding the performance degradation seen in RECOMP variants and Refiner at high  $k$  values. Also, we measured the impact on the efficiency of the RAG pipeline with token counts and end-to-end latency, confirming that EXIT significantly reduces context from 4,497.1 tokens to 594.4 tokens (86.8% fewer) at  $k = 30$ , even improving the quality. Also, EXIT’s latency scales nearly linearly (0.48s to 2.71s) and is much faster than the abstraction methods and the uncompressed baseline. These results demonstrate that EXIT consistently delivers significant accuracy improvements with minimal inference costs, regardless of the number of documents, making it well-suited for tasks involving larger retrieval sets.

## 6.2 Understanding End-to-End Latency Factors

While previous work has primarily focused on minimizing token counts to reduce reading time, we emphasize the importance of considering end-to-end latency, including compression, for building an efficient RAG pipeline. We provide a breakdown of the total end-to-end latency into reading time and compression time, along with an analysis of the average number of tokens in compressed documents, as shown in Figure 4. Although some methods achieve extreme token reduction—e.g., RECOMP-Abst (6.4%) and CompAct (10.2%)—their lengthy compression stages (1.46s and 7.99s, respectively) negate these gains, resulting in overall inference times that exceed the uncompressed baseline. By contrast, EXIT retains a moderate token ratio (28.3%) but completes com-

Table 2: Ablation studies on HQA examining (1) training data composition (Pos, H-Neg, Neg), (2) adaptive vs. fixed-length sentence selection, and (3) the impact of incorporating passage context during classification.

Configuration	EM $\uparrow$	F1 $\uparrow$	# token $\downarrow$
Ours (Pos + H-Neg + Neg)	<b>31.6</b>	<b>42.6</b>	195.1
Pos + H-Neg	30.0	41.3	286.8
Pos + Neg	29.8	40.9	404.6
w/o Adaptive Sentence Selection (2 sents)	29.4	40.7	<b>91.0</b>
w/o Adaptive Sentence Selection (4 sents)	30.2	41.4	166.5
w/o Passage as context	30.4	42.3	157.4

pression rapidly (0.36s), bringing total latency to 0.79s, a 23.3% improvement over the 1.03s needed without compression. This analysis highlights the importance of balancing token reduction and compression inference latency to achieve actual efficiency in the RAG pipeline. Also, our proposed framework, EXIT, aligns closely with these objectives, offering a practical method to context compression for RAG.

## 6.3 Ablation Studies

To better understand how the design choices in EXIT affect its overall performance and efficiency, we conduct ablation studies focusing on three key components: data sampling strategy, adaptive sentence selection, and context-aware extraction. Table 2 summarizes these results.

**Data Sampling Strategy.** Our training data combines positive, hard negative, and random negative samples to mirror the diversity of real-world retrieval scenarios. Compared to using only subsets of these sample types, the comprehensive strategy improves EM by 1.6 points over using only hard negatives and by 1.8 points over only random negatives. Relying solely on positive and random negatives led to excessive token retention, while depending only on hard negatives diminished the model’s ability to filter out spurious retrieval noise.

**Adaptive Sentence Selection.** We compare our adaptive selection mechanism to a fixed-length strategy, which limits selection to four sentences. While the fixed-length method reduces token count to 166.5, it lowers EM by 1.4 and F1 by 1.2 points. This highlights the importance of adapting sentence selection based on query complexity and document relevance rather than imposing a static constraint.

**Context-Aware Extraction.** We evaluate the effect of considering full document context when determining sentence relevance. Removing surrounding context reduces token usage by 38 tokens but decreases EM by 1.2 points, highlighting the importance of contextual awareness in preserving answer

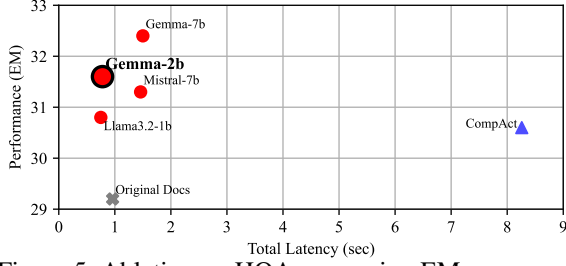


Figure 5: Ablation on HQA comparing EM scores and latency for different model configurations within EXIT (red dot). CompAct and Original Docs are included as indicators. Experiments on a single A100-80GB GPU.

accuracy. While context-aware extraction slightly increases token count, it ensures more precise and reliable responses.

These findings confirm that a balanced training data strategy enhances robustness, that adaptive sentence selection ensures efficiency, and that full document consideration preserves accuracy.

#### 6.4 Compressor Model Size and Strategy Impact

**Model Size Considerations.** Figure 5 presents an ablation study examining how different base models influence EM scores and total latency. Note that all models trained within EXIT achieve superior accuracy compared to uncompressed baselines and fast compression under 2 seconds. Specifically, Gemma-2B, our base classifier model, achieves a favorable balance between effectiveness and efficiency, delivering 31.6 EM points in 0.78s. When moving to the largest model, Gemma-7B, the highest accuracy (32.4 EM) is achieved, yet it also inflates latency to 1.50s, slightly exceeding the time of uncompressed documents. These results suggest that scaling up model parameters can improve performance but may also compromise latency benefits, emphasizing the flexibility of our proposed method by selecting an appropriate model depending on the user requirement.

**Abstractive vs. Extractive Compression.** Figure 5 also compares our extractive approach, EXIT, against CompAct, a 7B-scale abstractive compressor. Using Mistral-7B as the base model for both methods, EXIT (31.3 EM, 1.46s) significantly outperforms CompAct (30.6 EM, 8.26s) in terms of latency and maintains competitive accuracy. This stark difference underscores that the compression strategy, not the just model size, heavily influences efficiency. By relying on extraction rather than iterative summarization, EXIT capitalizes on a large-scale model to preserve high accuracy without incurring prohibitively long inference time.

Table 3: Performance comparison across different retrieval and reranking methods on 2WIKI.

Method	EM	F1	#tok.
<b>Passage-Level (Top-1)</b>			
bge-m3	14.6	22.7	156.5
bge-reranker-v2-gemma	15.0	22.9	156.2
RepLLaMA	16.0	24.1	157.2
UPR	9.2	16.8	155.3
Yes/No Ranking Prompt	13.4	21.8	157.5
Pairwise Ranking Prompt	13.2	21.5	156.2
<b>Sentence-Level (Top-5)</b>			
bge-m3	16.6	24.2	216.9
bge-reranker-v2-gemma	18.4	26.2	203.2
RepLLaMA	14.8	23.2	239.0
UPR	7.4	14.6	148.9
Yes/No ranking prompt	14.6	21.8	135.6
Pairwise ranking prompt	11.6	19.4	209.9
Ours	<b>24.8</b>	<b>30.1</b>	<b>150.2</b>

#### 6.5 Comparison with Reranking Methods

We compare EXIT with post-retrieval reranking methods in Table 3, including bge-m3 (Chen et al., 2024a), bge-reranker-v2-gemma (Li et al., 2023a), RepLLaMA (Ma et al., 2024), UPR (Sachan et al., 2022), Yes/No ranking prompting (Liang et al., 2023), and Pairwise ranking prompting (Qin et al., 2024). EXIT consistently achieves higher EM and F1 scores, demonstrating its ability to adaptively select query-relevant sentences while filtering out redundancy. By contrast, Top- $k$  re-ranking methods rely on static ranking, often retaining irrelevant content, increasing token count, and diminishing output quality. Moreover, they lack sentence-level context awareness, treating sentences independently without considering their broader document context, which can lead to incoherent or incomplete retrieval. These results underscore EXIT’s advantage in providing a more efficient and contextually relevant input for RAG than other post-retrieval methods.

#### 7 Conclusion

We present EXIT, an efficient context compression framework for RAG systems that leverages parallel processing and context-aware, adaptive sentence selection. Our experiments demonstrate that EXIT achieves superior performance across both single-hop and multi-hop QA tasks while maintaining practical inference speeds. Despite being trained only on HQA, EXIT shows a strong zero-shot generalization ability and proves effective across a wide range of open-source models of varying sizes. These results show that parallel extraction with smaller models outperforms larger abstractive methods, making it practical for real-world RAG.



## Limitation

This work focuses on RAG pipelines, and its effectiveness in general long-context scenarios like LongBench (Bai et al., 2024) remains future work. Our current approach relies on explicit sentence-level annotations to train the classifier. While these annotations were obtained manually in our experiments, they could potentially be automated through alternative means, such as by GPT-4 supervision or signals derived from the reader itself. The analysis in Section C.4 suggests that EXIT can effectively leverage such alternative supervision sources, further demonstrating its adaptability. However, we have not yet explored fully automated annotation strategies, which remains a promising direction for future research. Additionally, our study primarily focuses on a general-domain setting, leaving questions about the classifier’s performance in specialized domains unanswered. Investigating how well our approach generalizes to domain-specific or highly specialized corpora presents another valuable direction for future research. Lastly, we focus on a single-step RAG pipeline, where retrieval occurs only once, excluding more complex pipelines (Shao et al., 2023; Trivedi et al., 2023; Khattab et al., 2022). However, our proposed framework, EXIT, is orthogonal to these approaches and can be seamlessly integrated by compressing the retrieved documents from each retrieval step.

## Ethics Statement

This work enhances RAG-based QA without generating new content beyond what is retrieved. However, biases and inaccuracies in the source documents can still propagate through our compression process. Ensuring the reliability, fairness, and proper curation of underlying corpora is essential for ethical deployment. Future efforts should integrate bias detection, provenance tracking, and user-centric evaluations to promote more transparent and equitable real-world applications.

## 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, multi-task benchmark for long context understanding](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2024, Bangkok, Thailand, Au-

gust 11-16, 2024, pages 3119–3137. Association for Computational Linguistics.

Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024a. [BGE m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation](#). *CoRR*, abs/2402.03216.

Tong Chen, Hongwei Wang, Sihao Chen, Wenhao Yu, Kaixin Ma, Xinran Zhao, Hongming Zhang, and Dong Yu. 2024b. [Dense X retrieval: What retrieval granularity should we use?](#) In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 15159–15177. Association for Computational Linguistics.

Jianpeng Cheng and Mirella Lapata. 2016. [Neural summarization by extracting sentences and words](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics.

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](#). In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.

Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. 2023. [Adapting language models to compress contexts](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 3829–3846. Association for Computational Linguistics.

Eunsol Choi, Jennimaria Palomaki, Matthew Lamm, Tom Kwiatkowski, Dipanjan Das, and Michael Collins. 2021. [Decontextualization: Making sentences stand-alone](#). *Trans. Assoc. Comput. Linguistics*, 9:447–461.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic,

722	Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. <a href="#">The llama 3 herd of models</a> . <i>CoRR</i> , abs/2407.21783.	
738	Tao Ge, Jing Hu, Lei Wang, Xun Wang, Si-Qing Chen, and Furu Wei. 2024. <a href="#">In-context autoencoder for context compression in a large language model</a> . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
744	Michael Günther, Isabelle Mohr, Bo Wang, and Han Xiao. 2024. <a href="#">Late chunking: Contextual chunk embeddings using long-context embedding models</a> . <i>CoRR</i> , abs/2409.04701.	
748	Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. <a href="#">Constructing A multi-hop QA dataset for comprehensive evaluation of reasoning steps</a> . In <i>Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020</i> , pages 6609–6625. International Committee on Computational Linguistics.	
756	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. <a href="#">Lora: Low-rank adaptation of large language models</a> . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	
762	TaeHo Hwang, Soyeong Jeong, Sukmin Cho, SeungYoon Han, and Jong Park. 2024. <a href="#">DSLRL: Document refinement with sentence-level re-ranking and reconstruction to enhance retrieval-augmented generation</a> . In <i>Proceedings of the 3rd Workshop on Knowledge Augmented Methods for NLP</i> , pages 73–92, Bangkok, Thailand. Association for Computational Linguistics.	
769	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. <a href="#">Unsupervised dense information retrieval with contrastive learning</a> . <i>Trans. Mach. Learn. Res.</i> , 2022.	
774	Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong Park. 2024. <a href="#">Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity</a> . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , NAACL 2024, Mexico City, Mexico, June 16-21, 2024, pages 7036–7050. Association for Computational Linguistics.	780 781 782 783
	Soyeong Jeong, Jinheon Baek, Sung Ju Hwang, and Jong Park. 2023. <a href="#">Phrase retrieval for open domain conversational question answering with conversational dependency modeling via contrastive learning</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 6019–6031. Association for Computational Linguistics.	784 785 786 787 788 789 790 791
	Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023. <a href="#">Llmlingua: Compressing prompts for accelerated inference of large language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 13358–13376. Association for Computational Linguistics.	792 793 794 795 796 797 798 799
	Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2024. <a href="#">Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 1658–1677. Association for Computational Linguistics.	800 801 802 803 804 805 806 807 808
	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. <a href="#">Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension</a> . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers</i> , pages 1601–1611. Association for Computational Linguistics.	809 810 811 812 813 814 815 816
	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick S. H. Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. <a href="#">Dense passage retrieval for open-domain question answering</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 6769–6781. Association for Computational Linguistics.	817 818 819 820 821 822 823 824
	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. <a href="#">Generalization through memorization: Nearest neighbor language models</a> . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	825 826 827 828 829 830
	Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. <a href="#">Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive NLP</a> . <i>CoRR</i> , abs/2212.14024.	831 832 833 834 835
	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris Alberti,	836 837

838	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natural questions: a benchmark for question answering research</a> . <i>Trans. Assoc. Comput. Linguistics</i> , 7:452–466.	897
839		898
840		899
841		900
842		901
843		902
844		903
845	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. <a href="#">Efficient memory management for large language model serving with pagedattention</a> . In <i>Proceedings of the 29th Symposium on Operating Systems Principles, SOSP 2023, Koblenz, Germany, October 23-26, 2023</i> , pages 611–626. ACM.	904
846		
847		905
848		906
849		907
850		908
851		909
852		910
853	Jinhyuk Lee, Mujeen Sung, Jaewoo Kang, and Danqi Chen. 2021. <a href="#">Learning dense representations of phrases at scale</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021</i> , pages 6634–6647. Association for Computational Linguistics.	911
854		
855		912
856		913
857		914
858		915
859		916
860		917
861		918
862	Yuho Lee, Taewon Yun, Jason Cai, Hang Su, and Hwanjun Song. 2024. <a href="#">Unisumeval: Towards unified, fine-grained, multi-dimensional summarization evaluation for llms</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024</i> , pages 3941–3960. Association for Computational Linguistics.	
863		
864		919
865		920
866		921
867		922
868		
869	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. <a href="#">BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 7871–7880. Association for Computational Linguistics.	
870		
871		923
872		924
873		925
874		926
875		927
876		928
877		929
878	Patrick S. H. 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. 2020b. <a href="#">Retrieval-augmented generation for knowledge-intensive NLP tasks</a> . In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	930
879		931
880		932
881		933
882		934
883		935
884		936
885		937
886		938
887	Chaofan Li, Zheng Liu, Shitao Xiao, and Yingxia Shao. 2023a. <a href="#">Making large language models A better foundation for dense retrieval</a> . <i>CoRR</i> , abs/2312.15503.	939
888		940
889		
890	Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix X. Yu, and Sanjiv Kumar. 2023b. <a href="#">Large language models with controllable working memory</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 1774–1793. Association for Computational Linguistics.	941
891		942
892		943
893		944
894		945
895		
896		946
		947
		948
		949
		950
		951
		952
		953
		954
	Junyi Li, Xiaoxue Cheng, Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023c. <a href="#">Halueval: A large-scale hallucination evaluation benchmark for large language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 6449–6464. Association for Computational Linguistics.	
	Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. 2023d. <a href="#">Compressing context to enhance inference efficiency of large language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 6342–6353. Association for Computational Linguistics.	
	Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. 2023e. <a href="#">Compressing context to enhance inference efficiency of large language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 6342–6353. Association for Computational Linguistics.	
	Zhonghao Li, Xuming Hu, Aiwei Liu, Kening Zheng, Sirui Huang, and Hui Xiong. 2024. <a href="#">Refiner: Restructure retrieval content efficiently to advance question-answering capabilities</a> . <i>CoRR</i> , abs/2406.11357.	
	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soyulu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel J. Orr, Lucia Zheng, Mert Yuksekgönül, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2023. <a href="#">Holistic evaluation of language models</a> . <i>Trans. Mach. Learn. Res.</i> , 2023.	
	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. <a href="#">Lost in the middle: How language models use long contexts</a> . <i>Trans. Assoc. Comput. Linguistics</i> , 12:157–173.	
	Kun Luo, Zheng Liu, Shitao Xiao, Tong Zhou, Yubo Chen, Jun Zhao, and Kang Liu. 2024. <a href="#">Landmark embedding: A chunking-free embedding method for retrieval augmented long-context large language models</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 3268–3281. Association for Computational Linguistics.	



955	Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. 2024. <a href="#">Fine-tuning llama for multi-stage text retrieval</a> . In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2024, Washington DC, USA, July 14-18, 2024</i> , pages 2421–2425. ACM.	1013
956		1014
957		1015
958		1016
959		1017
960		1018
961		
962	Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Cristian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, and et al. 2024. <a href="#">Gemma: Open models based on gemini research and technology</a> . <i>CoRR</i> , abs/2403.08295.	1019
963		1020
964		1021
965		1022
966		1023
967		1024
968		
969		1025
970		1026
971		1027
972		1028
973		1029
974		1030
975		
976		1031
977		1032
978		1033
979		1034
		1035
980	Jesse Mu, Xiang Li, and Noah D. Goodman. 2023. <a href="#">Learning to compress prompts with gist tokens</a> . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1036
981		1037
982		1038
983		1039
984		
985		1040
986	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. <a href="#">Ranking sentences for extractive summarization with reinforcement learning</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)</i> , pages 1747–1759. Association for Computational Linguistics.	1041
987		1042
988		1043
989		1044
990		
991		1045
992		1046
993		1047
994		
995	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. <a href="#">MS MARCO: A human generated machine reading comprehension dataset</a> . In <i>Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, December 9, 2016</i> , volume 1773 of <i>CEUR Workshop Proceedings</i> . CEUR-WS.org.	1048
996		1049
997		1050
998		1051
999		1052
1000		1053
1001		1054
1002		1055
1003		1056
1004		
1005	Rodrigo Frassetto Nogueira and Kyunghyun Cho. 2019. <a href="#">Passage re-ranking with BERT</a> . <i>arXiv</i> , 1901.04085, abs/1901.04085.	1057
1006		1058
1007		1059
1008	OpenAI. 2023. <a href="#">GPT-4 technical report</a> . <i>arXiv</i> :2303.08774, abs/2303.08774.	1060
1009		1061
1010		1062
1011		1063
1012		
		1064
		1065
		1066
		1067
		1068
		1069
		1070
		1071
		1072
		1073
		1074
		1075
		1076
		1077
		1078
		1079
		1080
		1081
		1082
		1083
		1084
		1085
		1086
		1087
		1088
		1089
		1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200



1070	August 2, 2019, Volume 1: Long Papers, pages 4430–	
1071	4441. Association for Computational Linguistics.	
1072	Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie	
1073	Huang, Nan Duan, and Weizhu Chen. 2023. En-	
1074	hancing retrieval-augmented large language models	
1075	with iterative retrieval-generation synergy. In <i>Find-</i>	
1076	<i>ings of the Association for Computational Linguistics:</i>	
1077	<i>EMNLP 2023, Singapore, December 6-10, 2023,</i>	
1078	pages 9248–9274. Association for Computational	
1079	Linguistics.	
1080	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan	
1081	Scales, David Dohan, Ed H. Chi, Nathanael Schärli,	
1082	and Denny Zhou. 2023. Large language models can	
1083	be easily distracted by irrelevant context. In <i>Internat-</i>	
1084	<i>ional Conference on Machine Learning, ICML 2023,</i>	
1085	<i>23-29 July 2023, Honolulu, Hawaii, USA, volume</i>	
1086	<i>202 of Proceedings of Machine Learning Research,</i>	
1087	pages 31210–31227. PMLR.	
1088	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-	
1089	joon Seo, Richard James, Mike Lewis, Luke Zettle-	
1090	moyer, and Wen-tau Yih. 2024. REPLUG: retrieval-	
1091	augmented black-box language models. In <i>Proceed-</i>	
1092	<i>ings of the 2024 Conference of the North American</i>	
1093	<i>Chapter of the Association for Computational Lin-</i>	
1094	<i>guistics: Human Language Technologies (Volume 1:</i>	
1095	<i>Long Papers), NAACL 2024, Mexico City, Mexico,</i>	
1096	<i>June 16-21, 2024, pages 8371–8384. Association for</i>	
1097	<i>Computational Linguistics.</i>	
1098	Hwanjun Song, Hang Su, Igor Shalyminov, Jason Cai,	
1099	and Saab Mansour. 2024a. Finesure: Fine-grained	
1100	summarization evaluation using llms. In <i>Proceed-</i>	
1101	<i>ings of the 62nd Annual Meeting of the Association</i>	
1102	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	
1103	<i>pers), ACL 2024, Bangkok, Thailand, August 11-16,</i>	
1104	<i>2024, pages 906–922. Association for Computational</i>	
1105	<i>Linguistics.</i>	
1106	Hwanjun Song, Taewon Yun, Yuho Lee, Gihun Lee,	
1107	Jason Cai, and Hang Su. 2024b. Learning to	
1108	summarize from llm-generated feedback. <i>CoRR,</i>	
1109	abs/2410.13116.	
1110	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot,	
1111	and Ashish Sabharwal. 2023. Interleaving retrieval	
1112	with chain-of-thought reasoning for knowledge-	
1113	intensive multi-step questions. In <i>Proceedings of</i>	
1114	<i>the 61st Annual Meeting of the Association for</i>	
1115	<i>Computational Linguistics (Volume 1: Long Papers),</i>	
1116	<i>ACL 2023, Toronto, Canada, July 9-14, 2023, pages</i>	
1117	<i>10014–10037. Association for Computational Lin-</i>	
1118	<i>guistics.</i>	
1119	Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md. Rizwan	
1120	Parvez, and Graham Neubig. 2023. Learning to	
1121	filter context for retrieval-augmented generation.	
1122	<i>arXiv.2311.08377, abs/2311.08377.</i>	
1123	David Wingate, Mohammad Shoeybi, and Taylor	
1124	Sorensen. 2022. Prompt compression and contrastive	
1125	conditioning for controllability and toxicity reduction	
1126	in language models. In <i>Findings of the Association</i>	
	<i>for Computational Linguistics: EMNLP 2022, Abu</i>	1127
	<i>Dhabi, United Arab Emirates, December 7-11, 2022,</i>	1128
	pages 5621–5634. Association for Computational	1129
	Linguistics.	1130
	Zhenyu Wu, Chao Shen, and Meng Jiang. 2024. In-	1131
	structing large language models to identify and ig-	1132
	nore irrelevant conditions. In <i>Proceedings of the</i>	1133
	<i>2024 Conference of the North American Chapter of</i>	1134
	<i>the Association for Computational Linguistics: Hu-</i>	1135
	<i>man Language Technologies (Volume 1: Long Pa-</i>	1136
	<i>pers), NAACL 2024, Mexico City, Mexico, June 16-21,</i>	1137
	<i>2024, pages 6799–6819. Association for Computa-</i>	1138
	<i>tional Linguistics.</i>	1139
	Heming Xia, Zhe Yang, Qingxiu Dong, Peiyi Wang,	1140
	Yongqi Li, Tao Ge, Tianyu Liu, Wenjie Li, and Zhi-	1141
	fang Sui. 2024a. Unlocking efficiency in large lan-	1142
	guage model inference: A comprehensive survey of	1143
	speculative decoding. In <i>Findings of the Association</i>	1144
	<i>for Computational Linguistics, ACL 2024, Bangkok,</i>	1145
	<i>Thailand and virtual meeting, August 11-16, 2024,</i>	1146
	pages 7655–7671. Association for Computational	1147
	Linguistics.	1148
	Peng Xia, Kangyu Zhu, Haoran Li, Hongtu Zhu, Yun	1149
	Li, Gang Li, Linjun Zhang, and Huaxiu Yao. 2024b.	1150
	RULE: reliable multimodal RAG for factuality in	1151
	medical vision language models. In <i>Proceedings</i>	1152
	<i>of the 2024 Conference on Empirical Methods in</i>	1153
	<i>Natural Language Processing, EMNLP 2024, Miami,</i>	1154
	<i>FL, USA, November 12-16, 2024, pages 1081–1093.</i>	1155
	Association for Computational Linguistics.	1156
	Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2024. RE-	1157
	COMP: Improving retrieval-augmented LMs with	1158
	context compression and selective augmentation. In	1159
	<i>The Twelfth International Conference on Learning</i>	1160
	<i>Representations.</i>	1161
	Ruochen Xu, Song Wang, Yang Liu, Shuohang Wang,	1162
	Yichong Xu, Dan Iter, Pengcheng He, Chenguang	1163
	Zhu, and Michael Zeng. 2023. LMGQS: A large-	1164
	scale dataset for query-focused summarization. In	1165
	<i>Findings of the Association for Computational Lin-</i>	1166
	<i>guistics: EMNLP 2023, Singapore, December 6-10,</i>	1167
	<i>2023, pages 14764–14776. Association for Computa-</i>	1168
	<i>tional Linguistics.</i>	1169
	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-	1170
	gio, William W. Cohen, Ruslan Salakhutdinov, and	1171
	Christopher D. Manning. 2018. Hotpotqa: A dataset	1172
	for diverse, explainable multi-hop question answer-	1173
	ing. In <i>Proceedings of the 2018 Conference on Em-</i>	1174
	<i>pirical Methods in Natural Language Processing,</i>	1175
	<i>Brussels, Belgium, October 31 - November 4, 2018,</i>	1176
	pages 2369–2380. Association for Computational	1177
	Linguistics.	1178
	Chanwoong Yoon, Taewhoo Lee, Hyeon Hwang, Min-	1179
	byul Jeong, and Jaewoo Kang. 2024. Compact: Com-	1180
	pressing retrieved documents actively for question an-	1181
	swering. In <i>Proceedings of the 2024 Conference on</i>	1182
	<i>Empirical Methods in Natural Language Processing,</i>	1183
	<i>EMNLP 2024, Miami, FL, USA, November 12-16,</i>	1184

2024, pages 21424–21439. Association for Computational Linguistics.

Weijia Zhang, Jia-Hong Huang, Svitlana Vakulenko, Yumo Xu, Thilina Rajapakse, and Evangelos Kanoulas. 2024. [Beyond relevant documents: A knowledge-intensive approach for query-focused summarization using large language models](#). In *Pattern Recognition - 27th International Conference, ICPR 2024, Kolkata, India, December 1-5, 2024, Proceedings, Part XIX*, volume 15319 of *Lecture Notes in Computer Science*, pages 89–104. Springer.

Zheng Zhao, Shay B. Cohen, and Bonnie Webber. 2020. [Reducing quantity hallucinations in abstractive summarization](#). *CoRR*, abs/2009.13312.

Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and Jiawei Han. 2022. [Towards a unified multi-dimensional evaluator for text generation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 2023–2038. Association for Computational Linguistics.

## Appendix

In the Appendix, we provide additional implementation details and present supplementary results and analyses not covered in the main text.

### A More Related Work

#### A.1 Information Retrieval

Advancements in information retrieval (IR) have improved retrieval efficiency, granularity, and ranking quality. Encoder-based models, such as DPR (Karpukhin et al., 2020) and Contriever (Izacard et al., 2022), have been widely adopted for retrieval tasks. More recently, M3-Embedding (Chen et al., 2024a) unifies dense, multi-vector, and sparse retrieval via self-knowledge distillation, enabling retrieval across different text granularities.

Decoder-only LLMs have emerged as strong rerankers, leveraging large-scale data for improved ranking at the cost of computational complexity. LLaRA (Li et al., 2023a) and RePLlama (Ma et al., 2024) employ dense retrieval and multi-stage ranking, while UPR (Sachan et al., 2022) and ranking prompting methods (Liang et al., 2023; Qin et al., 2024) refine retrieval through generation-based approaches. However, these methods focus on reordering passages rather than compressing retrieved content.

Research on retrieval granularity (Seo et al., 2019; Lee et al., 2021; Jeong et al., 2023; Chen et al., 2024b; Hwang et al., 2024) explores balancing precision and contextual coherence. While fine-grained retrieval enhances specificity, excessive fragmentation can distort meaning (Choi et al., 2021). Retrieval and ranking aim to maximize recall but do not address post-retrieval redundancy. EXIT operates beyond retrieval and ranking, focusing on adaptive post-retrieval context compression, which selectively filters content to optimize both efficiency and relevance in RAG.

#### A.2 Summarization

Summarization methods are broadly categorized as extractive and abstractive. Extractive summarization retains factual consistency by selecting key sentences but often suffers from redundancy and coherence issues (Cheng and Lapata, 2016; Narayan et al., 2018). Abstractive summarization improves fluency and readability but is computationally intensive and prone to hallucination (See et al., 2017; Lewis et al., 2020a). Despite these challenges, its

ability to generate concise outputs has driven interest in its application (Song et al., 2024b; Lee et al., 2024; Song et al., 2024a).

Query-focused summarization refines document content based on relevance to a given query, removing unrelated information while preserving key details (Zhang et al., 2024; Xu et al., 2023). However, traditional summarization methods process single documents and prioritize completeness and faithfulness (Song et al., 2024a,b). In contrast, post-retrieval compression in RAG requires query-adaptive filtering across multiple retrieved documents. EXIT fills this gap by providing a dynamic and query-adaptive compression strategy, distinguishing between relevant, marginal, and irrelevant content while maintaining both latency efficiency and answer accuracy.

### B More Implementation Details

This section describes our training environment, data composition, and prompt templates. All experiments were conducted on an NVIDIA A100-SXM4-80GB GPU cluster. Training was performed on a single GPU with gradient accumulation.

#### B.1 Training Configuration

We employed LoRA (Hu et al., 2022) to fine-tune the compressor model, enabling parameter-efficient training while preserving performance. The model was trained with the following hyperparameters:

- Batch size: 8 per device
- Gradient accumulation steps: 8
- Learning rate: 1e-5
- Weight decay: 0.1
- Warmup ratio: 0.03
- Training epochs: 1
- Optimizer: `paged_adamw_8bit`
- Quantization: 4-bit with float16 precision
- LoRA configuration: Rank = 64, Scaling = 32, Dropout = 0.05

#### B.2 Model Selection and Training Time

Model selection was guided by validation loss. Training was conducted on our cluster, requiring approximately 90 hours.

#### B.3 Data Processing

We used SpaCy to segment documents into sentences. Table 4 shows the composition of the training and validation sets derived from HQA. The training set contains 427K sentences, including 213K positive (Pos), 107K hard-negative (H-Neg),

Table 4: Statistics of the training dataset constructed from HQA. Positive (Pos) sentences are required for the correct answer, Hard-Neg (H-Neg) sentences appear in the same passages but lack crucial evidence, and Neg sentences come from unrelated queries. Counts are in thousands (K).

Split	Pos	H-Neg	Neg	Total
Train	213K	107K	107K	427K
Valid	2.4K	1.2K	1.2K	4.8K

and 107K negative (Neg) instances. The validation set includes 4.8K sentences with a similar distribution. This balanced composition ensures the classifier encounters diverse retrieval scenarios during training.

#### B.4 Inference Settings

For inference, we set:

- Temperature: 0.0
- Top-p: 1.0
- vLLM version: v0.5.5
- Relevance threshold ( $\tau$ ): 0.5

#### B.5 Prompt Templates

Table 5: A prompt template for document compression.

Compression Prompt Template
Query: {query} Full context: {original passage} Sentence: {sentence} Is this sentence useful in answering the query? Answer only “Yes” or “No”.

Full prompt templates for compression and QA tasks are provided in Tables 5 and 6, respectively.

#### B.6 Reproducibility

We will release our codebase, including dataset pre-processing scripts, evaluation protocols, and model checkpoints, upon publication. All random seeds are set to 42 to facilitate reproducibility.

Table 6: A prompt template used by the reader model for the QA task.

QA Prompt Template
Context information is below. <hr/> {context} <hr/> Given the context information and not prior knowledge, answer the query. Do not provide any explanation. Query: {query} Answer:

Table 7: Performance comparison between in-domain (HQA) and out-of-domain (2WIKI) datasets using BM25 as the retriever model. Best results are highlighted in **bold**, and second best results are underlined.

Compressor	Type	HQA			2WIKI		
		EM $\uparrow$	F1 $\uparrow$	#token (%) $\downarrow$	EM $\uparrow$	F1 $\uparrow$	#token (%) $\downarrow$
Top-5 Documents							
Original Docs	-	28.2	39.1	755.8 (100.0)	19.6	25.9	789.7 (100.0)
RECOMP-Abst	Abs.	27.8	39.2	<b>63.0 (8.3)</b>	<b>25.0</b>	<b>30.6</b>	<b>55.7 (7.1)</b>
CompAct	Abs.	30.6	41.2	76.9 (10.2)	19.6	<u>29.1</u>	71.0 (9.0)
Refiner	Abs.	<u>30.8</u>	<u>42.3</u>	84.0 (11.1)	19.6	27.8	62.9 (8.0)
RECOMP-Extr	Ext.	28.2	37.5	93.1 (12.3)	11.8	19.1	99.1 (12.5)
LongLLMLingua	Ext.	28.6	39.7	230.3 (30.5)	22.2	27.4	236.6 (30.0)
EXIT (Ours)	Ext.	<b>33.4</b>	<b>44.5</b>	177.8 (23.5)	<u>24.4</u>	<u>29.1</u>	138.2 (17.5)
Top-20 Documents							
Original Docs	-	31.2	42.5	3009.5 (100.0)	23.0	<u>30.6</u>	3132.5 (100.0)
RECOMP-Abst	Abs.	29.0	39.7	<b>69.7 (2.3)</b>	22.8	28.2	<b>52.3 (1.7)</b>
CompAct	Abs.	<u>31.6</u>	<u>43.0</u>	109.5 (3.6)	20.4	28.7	113.2 (3.6)
Refiner	Abs.	28.4	38.8	136.1 (4.5)	18.8	26.9	108.4 (3.5)
RECOMP-Extr	Ext.	28.4	37.0	93.5 (3.1)	11.2	19.3	96.7 (3.1)
LongLLMLingua	Ext.	31.0	40.7	558.0 (18.5)	<u>23.6</u>	28.3	593.7 (19.0)
EXIT (Ours)	Ext.	<b>35.2</b>	<b>46.9</b>	411.3 (13.7)	<b>27.2</b>	<b>32.4</b>	312.7 (10.0)

## C Additional Experimental Results and Analyses

Table 15 shows detailed results for each dataset and model configuration under zero-shot QA prompts, using both Top-5 and Top-20 retrieval. For the few-shot QA prompts, Table 16 summarizes the results, where we randomly selected five training examples per dataset as demonstrations. Table 17 provides comprehensive token count statistics, and Table 18 breaks down end-to-end latency.

### C.1 Performance with Sparse Retrieval

To assess EXIT’s robustness with different retrieval architectures, we evaluate it with BM25, a sparse retrieval method. Table 7 compares EXIT’s performance on HQA (in-domain) and 2WIKI (out-of-domain) under Top-5 and Top-20 retrieval settings.



Table 8: Performance comparison between in-domain (HQA) and out-of-domain (2WIKI) datasets using GPT-4o as the reader model. Best results are highlighted in **bold**, and second best results are underlined.

Compressor	Type	HQA			2WIKI		
		EM $\uparrow$	F1 $\uparrow$	#token (%) $\downarrow$	EM $\uparrow$	F1 $\uparrow$	#token (%) $\downarrow$
Top-5 Documents							
Original Docs	-	37.2	<u>48.6</u>	735.3 (100.0)	<u>31.2</u>	<u>35.3</u>	764.5 (100.0)
RECOMP-Abst	Abs.	29.4	40.2	<b>62.8 (8.5)</b>	23.8	27.8	<b>53.5 (7.0)</b>
CompAct	Abs.	<u>37.4</u>	48.0	74.3 (10.1)	30.0	33.6	67.6 (8.8)
Refiner	Abs.	35.8	47.5	<u>71.4 (9.7)</u>	27.8	32.6	<u>54.2 (7.1)</u>
RECOMP-Extr	Ext.	32.4	41.7	87.1 (11.8)	25.6	28.6	93.6 (12.2)
LongLLMLingua	Ext.	34.2	45.2	223.1 (30.3)	30.6	34.5	230.5 (30.2)
EXIT (Ours)	Ext.	<b>38.2</b>	<b>50.4</b>	191.2 (26.0)	<b>31.8</b>	<b>35.8</b>	145.6 (19.0)
Top-20 Documents							
Original Docs	-	<b>39.6</b>	<b>51.8</b>	2940.5 (100.0)	<b>40.0</b>	<b>43.8</b>	3066.2 (100.0)
RECOMP-Abst	Abs.	33.6	44.2	<b>62.7 (2.1)</b>	26.6	32.1	<b>48.9 (1.6)</b>
CompAct	Abs.	33.0	43.7	106.0 (3.6)	23.0	27.3	105.1 (3.4)
Refiner	Abs.	31.8	41.5	130.6 (4.4)	31.0	35.5	100.3 (3.3)
RECOMP-Extr	Ext.	31.2	39.6	<u>86.2 (2.9)</u>	23.2	27.2	<u>91.0 (3.0)</u>
LongLLMLingua	Ext.	38.8	49.4	549.5 (18.7)	35.4	39.6	581.5 (19.0)
EXIT (Ours)	Ext.	<u>39.4</u>	<u>50.1</u>	453.6 (15.4)	<u>35.6</u>	<u>40.3</u>	346.7 (11.3)

With Top-5 retrieval, EXIT shows notable gains on HQA, improving EM (33.4 vs. 28.2) and F1 (44.5 vs. 39.1) over the uncompressed baseline. Although RECOMP-Abst performs best on 2WIKI (25.0 EM, 30.6 F1), EXIT remains competitive (24.4 EM, 29.1 F1).

EXIT’s advantages grow with Top-20 retrieval. On HQA, EXIT outperforms all baselines, improving EM by 4.0 points (35.2 vs. 31.2) and F1 by 4.4 points (46.9 vs. 42.5) compared to using uncompressed documents. On 2WIKI, EXIT achieves the highest scores (27.2 EM, 32.4 F1), confirming its generalizability across domains and retrieval strategies.

## C.2 Performance with Proprietary Model

We further examine EXIT’s effectiveness using GPT-4o as the reader. Table 8 compares performance on HQA (in-domain) and 2WIKI (out-of-domain), along with compression rates.

For Top-5 retrieval, EXIT attains the best accuracy on HQA (38.2 EM, 50.4 F1) while retaining only 26.0% of tokens. This surpasses the uncompressed baseline (37.2 EM, 48.6 F1) with a 74% token reduction. On 2WIKI, EXIT maintains leading accuracy (31.8 EM, 35.8 F1) while using just 19.0% of the original tokens.

Under Top-20 retrieval, where uncompressed documents benefit from greater coverage, EXIT still achieves competitive accuracy with substantially fewer tokens. On HQA, EXIT closely matches the uncompressed EM score (39.4 vs. 39.6) while using only 15.4% of tokens. Although

Table 9: Comparison of EXIT with baselines on QA tasks on HQA. EXIT (GPT-4o) refers to a variant using GPT-4o as a zero-shot relevance classifier, while EXIT (Ours) represents our proposed method trained with supervised data.

Method	EM	F1	#token
Original Docs	29.2	40.2	735.2
RECOMP-Extr	25.2	34.8	89.7
RECOMP-Abst	19.8	24.7	55.5
LongLLMLingua	27.4	40.2	227.1
CompAct	29.4	40.9	76.3
Refiner	28.8	40.7	73.4
EXIT (GPT-4o)	30.4	41.9	84.2
EXIT (Ours)	<b>31.6</b>	<b>42.6</b>	191.0

RECOMP variants compress more aggressively, they suffer marked performance drops. LongLLMLingua performs similarly to EXIT but retains more tokens (18.7% vs. 15.4%).

These findings illustrate EXIT’s ability to balance performance and efficiency, making it valuable for API-based proprietary models where token costs and accuracy both matter.

## C.3 Impact of Threshold $\tau$

We analyze EXIT’s sensitivity to the relevance threshold  $\tau$ . Figure 6 shows EXIT’s performance across various  $\tau$  values.

EXIT remains stable over a wide threshold range, with strong results between  $\tau=0.3$ – $0.5$ . At  $\tau=0.3$ , EXIT reaches 25.2 EM using only 25% of the tokens (195.82 vs. 780.95 for the baseline), a substantial improvement over the original documents (18.0 EM). F1 scores also remain consistently higher than the baseline (30.21–30.73 vs. 25.74).

Even under extreme compression ( $\tau=0.9$ , 7.9% of tokens), EXIT achieves better accuracy (24.0 EM, 29.44 F1) than the uncompressed documents. Conversely, a lenient threshold ( $\tau=0.1$ ) retains more tokens but still provides benefits, demonstrating that EXIT effectively identifies crucial content under varying conditions.

This robustness across thresholds gives practitioners flexibility to adjust the compression-accuracy trade-off without severely impacting performance.

## C.4 Adaptability of EXIT to Different Supervision Signals

A key strength of EXIT is its ability to integrate different classifiers for sentence selection without being constrained to manually annotated datasets.

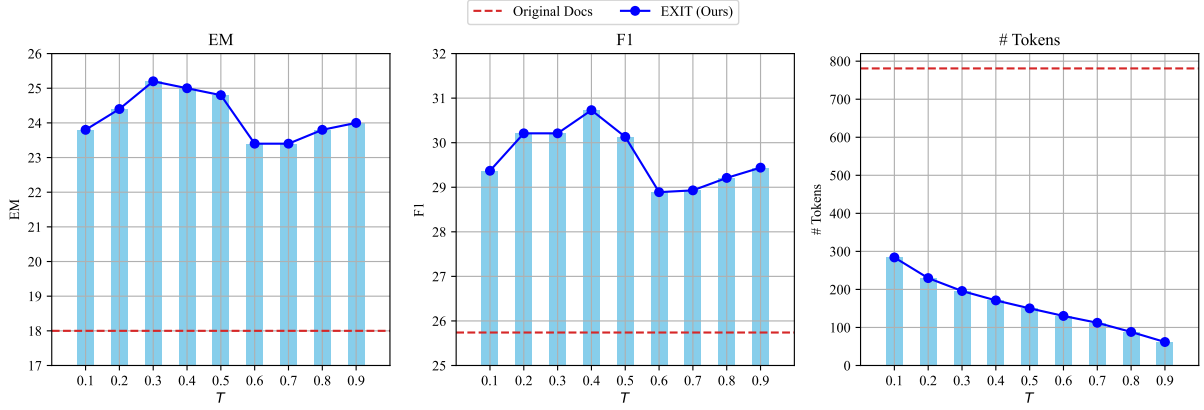


Figure 6: Changes in EM, F1 score, and token count as the threshold  $\tau$  for retaining sentences is adjusted.

Table 10: Classification performance for Yes/No labels.

Overall			
Class	Precision $\uparrow$	Recall $\uparrow$	F1-Score $\uparrow$
Yes	0.91	<b>0.93</b>	<b>0.92</b>
No	<b>0.93</b>	0.91	<b>0.92</b>
Hard Negative			
Yes	0.86	<b>0.93</b>	<b>0.89</b>
No	<b>0.93</b>	0.84	0.88
Negative			
Yes	<b>0.96</b>	0.93	<b>0.95</b>
No	0.93	<b>0.96</b>	<b>0.95</b>

		HQA (In-Domain)		2WIKI (Out-of-Domain)	
Actual	Yes	0.93	0.07	0.76	0.24
	No	0.09	0.91	0.06	0.94
		Predicted		Predicted	
		Yes	No	Yes	No

Figure 7: Confusion matrices (row-normalized) for context-aware relevance classification on HQA (in-domain) and 2WIKI (out-of-domain).

To explore this flexibility, we evaluate an alternative approach where relevance scores are derived from GPT-4o without explicit fine-tuning. As shown in Table 9, EXIT maintains strong performance, outperforming or closely matching baselines such as LongLLMLingua and CompAct. This demonstrates that EXIT can leverage diverse supervision signals while remaining adaptable to different scoring mechanisms. Moreover, these results suggest the potential of utilizing large-scale pseudo-labeled data to further refine EXIT’s training, enhancing scalability without relying strictly on human-labeled datasets. This adaptability highlights EXIT’s robustness and practical applicability across various retrieval settings.

### C.5 Analysis of Classification Performance Across Negative Sample Types

Table 10 presents EXIT’s sentence-level classification performance, broken down by negative sample type. EXIT achieves 0.92 F1 for both positive (“Yes”) and negative (“No”) classes overall, indicating a balanced ability to identify essential and non-essential sentences.

The model excels at filtering random negatives (0.95 F1), effectively discarding irrelevant content. With hard negatives (topically related but not essential), EXIT still performs well (0.89 F1 for positive, 0.88 for negative), handling nuanced relevance distinctions.

These results highlight EXIT’s adaptability and confirm its suitability for real-world RAG scenarios where both overtly irrelevant and subtly extraneous content must be managed.

### C.6 Classification Performance

To better understand the effectiveness of our context-aware relevance classifier, we report row-normalized confusion matrices for both in-domain (HQA) and out-of-domain (2WIKI) datasets, as shown in Figure 7. On HQA, the classifier displays a perfectly balanced ability to recognize both relevant (“Yes”) and irrelevant (“No”) sentences, achieving over 90% precision and recall in each category. While, on the 2WIKI dataset, the classifier exhibits a slight drop in recall for “Yes” sentences, it still performs strong classification ability with over 70% recall and 90% precision.

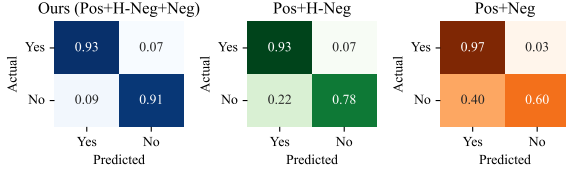


Figure 8: Row-normalized confusion matrices for classification performance under different training data conditions: Ours (Pos+H-Neg+Neg), Pos+H-Neg, and Pos+Neg. Each matrix compares the predicted (“Yes”/“No”) labels against the actual labels.

These results confirm that the classifier performs robustly in its training domain and generalizes reasonably well to unseen queries, yet we leave narrowing this discrepancy as a valuable future research direction.

## C.7 Classification Performance under Ablation Settings

### C.7.1 Analysis of Training Data Composition

Figure 8 presents row-normalized confusion matrices comparing classification performance across three training data configurations: Ours (Pos+H-Neg+Neg), Pos+H-Neg, and Pos+Neg. Under the Ours setup, the classifier displays a balanced ability to identify both “Yes” (relevant) and “No” (irrelevant) sentences, achieving an F1-score of 0.92 for both classes. In contrast, excluding one type of negative sample (Pos+H-Neg or Pos+Neg) reduces overall robustness, evidenced by declines in both accuracy and class-wise F1 scores. For instance, the Pos+Neg configuration struggles to maintain balance, accurately identifying “Yes” instances but misclassifying a substantial number of “No” cases. These results confirm that incorporating a comprehensive mix of positive, hard-negative, and random-negative samples leads to more reliable and contextually aware sentence selection, thereby improving the classifier’s performance in practical retrieval-augmented QA scenarios.

### C.7.2 Impact of Context on Classification Performance

We evaluate the classifier’s performance with and without broader passage-level context. Figure 9 shows that including context maintains over 90% precision and recall for both “Yes” and “No” classes. Without context, precision and recall decline, weakening the distinction between relevant and irrelevant sentences. This emphasizes the importance of incorporating passage-level context for accurately identifying answer-critical information.

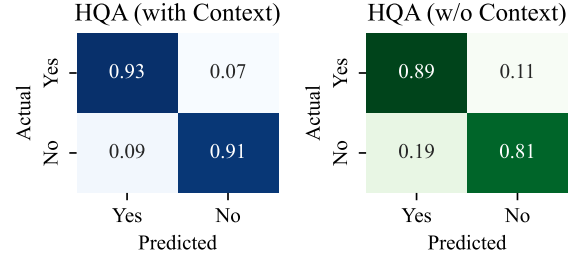


Figure 9: Row-normalized confusion matrices comparing classification performance with (left) and without (right) contextual information on HQA. The availability of context improves the model’s ability to accurately distinguish relevant (“Yes”) from non-relevant (“No”) sentences.

Table 11: Training data ablation comparison. Best results per metric are highlighted in **bold**.

Training Data	EM $\uparrow$	F1 $\uparrow$	# token $\downarrow$
HQA	<b>31.6</b>	<b>42.6</b>	195.1
2WIKI	29.2	40.3	<b>135.3</b>
2WIKI+HQA	30.6	42.0	232.2

## C.8 Training Data Ablation Analysis

Table 11 compares models trained on HQA, 2WIKI, or both. Training solely on HQA yields the highest EM and F1 (31.6 EM, 42.6 F1) with moderate token usage. In contrast, 2WIKI training improves compression but lowers accuracy (29.2 EM, 40.3 F1). Combining datasets does not surpass HQA alone.

This finding suggests that data quality and structure matter more than quantity. HQA’s annotations appear particularly effective for learning robust compression strategies that generalize well, validating our choice to use it as the primary training dataset.

## C.9 Case Studies

To illustrate how EXIT’s extractive compression strategy improves both accuracy and readability, we present qualitative examples and comparisons with other compression methods.

**Original Documents vs. Ours.** In Table 12, the original documents contain the correct answer (“Custard”) but also include distracting information (“Eggnog”). Despite having the necessary evidence, the reader fails to produce the correct answer, likely due to this distractor. In contrast, our method (Ours) filters out irrelevant details, drastically reduces input length, and retains only the essential context

needed to answer the query accurately. As a result, the reader confidently generates the correct answer (“Custard”).

In a second scenario (Table 13), the original documents retrieve multiple documents related to “Wagner” or “sci-fi” series but fail to provide any content explicitly linking James Belushi to the correct 90s sci-fi series, resulting in an incorrect prediction. Surprisingly, Ours removes all retrieved context entirely, providing the reader with no additional information. Under this no-context condition, the reader relies solely on its internal knowledge and, in this case, correctly identifies “Wild Palms.” While this outcome indicates a form of hallucination or model bias—since the answer emerges without external supporting evidence—it also demonstrates Ours’ capability to avoid misleading context. By eliminating irrelevant or confusing documents, Ours can sometimes allow the model’s internal knowledge to surface, leading to correct answers even in the absence of any retrieved information.

**Comparisons with Other Methods.** Table 14 compares Ours with several competing compression approaches. CompAct preserves some relevant information but introduces hallucinations, causing the reader to claim that it cannot find the correct answer. Refiner omits the crucial entity required to answer the query, demonstrating how abstractive compressors may inadvertently remove key content. In contrast, Ours avoids hallucinations and retains the answer’s entity in a concise, coherent form.

LongLLMLingua’s token-level filtering approach yields unreadable text and removes the essential “Romania” entity, preventing the reader from generating the correct answer. Ours, on the other hand, maintains semantic coherence and includes the correct entity, allowing the reader to produce the correct answer without interference.

These case studies highlight the advantages of Ours: it eliminates distractors, preserves critical entities, and maintains semantic integrity. Consequently, the reader consistently arrives at correct answers with fewer tokens and no hallucinations.



Table 12: Case study comparing compressed contexts and answers between Original Docs and Ours.

	Original Docs	Ours
<b>Query</b>	Crème Anglaise is the French version of which English dessert item?	
<b>Context</b>	<p>[1] Creme anglaise Crème anglaise Crème anglaise (French for ""English cream"" ) is a light pouring <b>Cus-tard</b> used as a dessert cream or sauce. It is a mix of sugar, egg yolks, oil, and hot milk</p> <p>...</p> <p>[2] Creme anglaise However, the ice cream base is much thicker and has various flavourings. The American South it is known as ""Custard."" It can be served like <b>Eggnog</b> during the Christmas season. Other names include the French terms ""crème à l'anglaise"" (""English-style cream"" ) and ""crème française"" (""French cream"" ). Crème anglaise Crème anglaise (French for ""English cream"" ) is a light pouring Custard used as a dessert cream or sauce. It is a mix of sugar, egg yolks, oil, and hot milk often flavoured with vanilla. Its name may derive from the prevalence of sweet Custards in English desserts. The cream is made by</p> <p>[3] Creme anglaise</p> <p>...</p> <p>[4] Custard lemon. ""Crème pâtissière"" is a key ingredient in many French desserts including mille-feuille (or Napoleons) and filled tarts. It is also used in Italian pastry and sometimes in Boston cream pie. The thickening of the Custard is caused by the combination of egg and cornstarch. Corn flour or flour thicken at 100 °C and as such many recipes instruct the pastry cream to be boiled. In a traditional Custard such as a ""crème anglaise"" , where egg is used alone as a thickener, boiling results in the over cooking and subsequent 'curdling' of the Custard; however, in a pastry cream, starch</p> <p>[5] Cremeschnitte usually pure thick Custard, less commonly combined with meringue (whipped egg whites and sugar) creme. A similar recipe with only meringue filling is called Šampita. In Australia, the dish is more commonly known as a 'vanilla slice'. Cremeschnitte A cremeschnitte (,,,,,, ) is a chantilly and Custard cream cake dessert popular in several Central European countries. There are many regional variations, but they all include a puff pastry base and Custard cream. In Slovenia, kremna rezina is commonly associated with the town of Bled, an Alpine tourist destination in northwestern Slovenia. The recipe</p>	<p>[1] Creme anglaise Crème anglaise Crème anglaise (French for ""English cream"" ) is a light pouring <b>Cus-tard</b> used as a dessert cream or sauce.</p> <p>[2] Creme anglaise Other names include the French terms ""crème à l'anglaise"" (""English-style cream"" ) and ""crème française"" (""French cream"" ). Crème anglaise Crème anglaise (French for ""English cream"" ) is a light pouring Custard used as a dessert cream or sauce.</p> <p>[3] Creme anglaise Alternatively, it can be drunk as a dessert on its own, for example in ""Île flottante"" (""floating island"" ): the cream is poured into a bowl with a piece of meringue (""blancs en neige"" ) floated on top along with praline. It can also be used as a base for desserts such as ice cream or crème brûlée.</p> <p>[4] Custard ""Crème pâtissière"" is a key ingredient in many French desserts including mille-feuille (or Napoleons) and filled tarts.</p>
<b>Answer</b>	Custard	
<b>Predict</b>	<b>Eggnog</b>	<b>Custard</b>

Table 13: Case study comparing compressed contexts and answers between Original Docs and Ours. Despite containing no relevant context, Ours method predicts the correct answer, indicating a hallucination scenario.

	Original Docs	Ours
<b>Query</b>	Which 90s sci fi series with James Belushi was based on Bruce Wagner's comic strip of the same name?	
<b>Context</b>	<p>[1] Michael I. Wagner patterned his character after Wagner's mannerisms and physical behavior. The series ran on Thursday nights in the Spring of 1988 during the same time slot as NBC's "The Cosby Show", and with that competition could not attract a sufficient audience to get renewed for the following season.</p> <p>...</p> <p>Wagner helped develop and write the Bochco animated series "Capitol Critters", he also wrote and served as supervising producer</p> <p>[2] Michael I. Wagner Steven Bochco and several of his projects. Wagner was asked by ABC in 1987 to help develop a new science fiction series, " <b>Probe</b> ", a light-hearted series about a scientific crime fighter named Austin James.</p> <p>...</p> <p>Parker Stevenson, who played the lead character, stated in a later interview that he</p> <p>[3] John Wagner the mid-1990s Wagner worked on a number of licensed properties for Dark Horse Comics in the US, including "Aliens", "Star Wars" – notably solo stories starring Boba Fett and the comics strand of the multimedia project "" – and "".</p> <p>...</p> <p>It was nominated for the Angoulême International Comics Festival Prize for Scenario in 2006. In 2000 Wagner</p> <p>[4] Martin Wagner (artist) Martin Wagner (artist) Martin Wagner (born April 27, 1966) is an American artist, cartoonist, and filmmaker.</p> <p>...</p> <p>. His production schedule became increasingly protracted and he ceased publishing the series altogether following issue No. 12 in 1994. In 1996 he made a</p> <p>[5] Wired (film) "L.A. Law", "Murphy Brown", and "Seinfeld"), Chiklis gained fame for portraying the lead roles of Commissioner Tony Scali on the ABC police drama "The Commish" (1991-1996), and LAPD Detective Vic Mackey on the FX police drama "The Shield" (2002-2008) and as well as Marvel superhero Ben "The Thing" Grimm in the films "Fantastic Four" (2005) and "" (2007).</p> <p>...</p> <p>Wired (film) Wired is a 1989 biographical film of comedian and actor John Belushi, directed by Larry Peerce, and adapted from the 1984 book of the same name by "Washington"</p>	No context (completely pruned)
<b>Answer</b>	Wild Palms	
<b>Predict</b>	<b>Probe</b>	Wild Palms

Table 14: Combined case studies comparing compressed contexts and answers across different methods.

Case Study 1	CompAct	Ours
<b>Query</b>	Where was there a series of explosions and fire on 23 October 1989 (killing 23 employees and injuring 314), an explosion on 23 June 1999 (killing two men and injuring three), and an explosion and fire on 27 March 2000 (causing 1 death and 71 injuries)?	
<b>Context</b>	The Flixborough disaster was an explosion at a chemical plant near Flixborough, England on 1 June 1974. The incident started on October 23, 1989, at 1400 Jefferson Road, Pasadena, Texas .	[1] Phillips disaster of 1989 On October 23, 1989, at 1400 Jefferson Road, Pasadena, Texas , a powerful explosion killed 23 and injured 314. A powerful and devastating explosion and fire ripped through the HCC, killing 23 persons—all working at the facility—and injuring 314 other (185 Phillips 66 employees and 129 contract employees).
<b>Answer</b>	Pasadena, Texas	
<b>Predict</b>	I couldn't find any information about the specific incidents you mentioned. The Flixborough disaster occurred on 1 June 1974.	Pasadena, Texas
Case Study 2	Refiner	Ours
<b>Query</b>	Who said that there are only three true sports, bullfighting, mountaineering, and motor racing, the rest are merely games?	
<b>Context</b>	1.1. Extreme sport There are only three sports: bullfighting, motor racing, and mountaineering; all the rest are merely games.	[1] Extreme sport usually, but wrongly, attributed to Ernest Hemingway . The phrase is: There are only three sports: bullfighting, motor racing, and mountaineering; all the rest are merely games.
<b>Answer</b>	Ernest Hemingway	
<b>Predict</b>	There is no information provided about who said that.	Ernest Hemingway
Case Study 3	LongLLMLingua	Ours
<b>Query</b>	Which was the only Eastern Bloc country to participate in the 1984 LA Olympics?	
<b>Context</b>	984 Summer Olympics the majority of Soviet Bloc will not participate Games, as expected. 984 Summer boycott The boycott of the 1984 Summer in Los Angeles years after U.S.-led boycott of the 1980 Summer Olympics in The boycott involved 14 Bloc countries and allies, led by the Soviet Union, which initiated boycott on May 8, 1984. Booting countries organized major event, the Friendship Games, in July and August 1984 Although the boycott by the Soviet Union affected number Olympic Summer were perhaps the first games of a new era to make a profit. Although a boycott led by the Soviet Union depleted the field in certain sports, 140 National Olympic Committees took part, which was a record at the time. Again, without the participation of the Eastern European countries, the 1984 Games were dominated by their host country. The Games were also the first time mainland China (People's Republic) participated. According to British journalist Andrew Jennings, a KGB colonel stated that the agency's officers had posed as anti-doping authorities from	[1] 1984 Summer Olympics boycott The boycott involved 14 Eastern Bloc countries and allies, led by the Soviet Union, which initiated the boycott on May 8, 1984. ... [3] Summer Olympic Games Eastern Bloc that did attend the 1984 Olympics. Although a boycott led by the Soviet Union depleted the field, 140 NOCs took part. Without Eastern European countries, the 1984 Games were dominated by the host. The Games were also the first time mainland China participated. ... [5] 1984 Summer Olympics boycott However, no threat to Eastern Bloc athletes was discovered, and the athletes from the Eastern Bloc country that did attend the 1984 games—Romania—encountered no problems.
<b>Answer</b>	Romania	
<b>Predict</b>	China	Romania

Table 15: Zero-shot QA prompt evaluation of compressor performance across different Top- $k$  scenarios, models, and datasets, measured by EM, F1, and inference latency (Lat.). 8B reader experiments were conducted on a single A100-80GB GPU, while 70B reader experiments utilized 4 A100-80GB GPUs in parallel. Best results for each dataset are highlighted in **bold**, Second best results are highlighted in underline.

Compressor	Type	NQ			TQA			HQA			2WIKI			AVG.		
		EM ↑	F1 ↑	Lat. ↓	EM ↑	F1 ↑	Lat. ↓	EM ↑	F1 ↑	Lat. ↓	EM ↑	F1 ↑	Lat. ↓	EM ↑	F1 ↑	Lat. ↓
Llama-3.1-8B-Instruct																
Top-5 Documents																
Original Docs	-	34.6	47.1	1.0	58.8	68.6	0.9	28.1	38.6	1.0	16.1	24.9	1.1	34.4	44.8	1.0
RECOMP-Abst	Abs.	31.3	43.2	1.6	55.9	65.7	1.4	26.5	37.0	2.2	22.7	29.1	2.1	34.1	43.7	1.8
CompAct	Abs.	32.9	44.6	8.5	58.1	67.7	8.8	28.8	39.8	8.3	16.8	26.0	8.1	34.2	44.5	8.4
Refiner	Abs.	32.9	45.0	28.1	59.2	68.9	10.9	28.8	40.0	6.9	16.8	25.4	6.4	34.4	44.8	13.1
RECOMP-Extr	Ext.	34.6	44.6	0.5	56.5	65.1	0.4	23.4	32.8	0.4	11.2	19.6	0.6	31.4	40.5	0.5
LongLLMLingua	Ext.	30.2	41.5	0.9	59.4	68.0	0.8	28.0	38	0.8	21.5	27.4	0.9	34.8	43.7	0.9
Ours (EXIT)	Ext.	35.9	47.8	0.8	60.8	69.9	0.7	30.6	41.5	0.8	24.2	30.8	0.9	37.9	47.5	0.8
Top-20 Documents																
Original Docs	-	36.6	49.5	3.4	62.0	71.7	2.9	29.9	40.5	2.9	18.8	27.9	3.2	36.8	47.4	3.1
RECOMP-Abst	Abs.	26.9	38.3	1.7	57.3	66.6	1.9	26.8	37.1	2.4	22.7	28.8	2.6	33.4	42.7	2.2
CompAct	Abs.	33.8	45.4	26.1	57.8	67.5	24.5	28.9	39.6	27.5	16.7	24.6	32.2	34.3	44.3	27.6
Refiner	Abs.	30.1	41.4	28.7	57.6	67.0	44.2	26.6	37.3	29.6	16.3	24.9	10.8	32.7	42.6	28.3
RECOMP-Extr	Ext.	32.8	42.6	0.6	55.5	63.6	0.4	22.2	31.2	0.5	10.0	18.3	0.7	30.1	38.9	0.6
LongLLMLingua	Ext.	33.4	45.1	2.8	62.4	71.2	2.7	31.2	41.4	2.8	24.1	30.1	2.9	37.8	46.9	2.8
Ours (EXIT)	Ext.	38.1	50.8	1.8	62.8	72.0	1.7	32.9	44.0	1.8	25.5	32.3	2.0	39.8	49.8	1.8
Llama-3.1-70B-Instruct																
Top-5 Documents																
Original Docs	-	35.6	48.0	8.6	65.1	73.9	7.7	33.7	44.5	8.3	20.8	28.3	9.1	38.8	48.7	8.4
RECOMP-Abst	Abs.	34.1	47.0	4.5	61.3	70.6	3.3	30.3	40.8	4.4	24.2	30.3	4.2	37.5	47.2	4.1
CompAct	Abs.	34.1	45.4	11.9	62.6	71.1	11.7	33.8	44.1	11.0	20.5	27.4	11.6	37.8	47.0	11.5
Refiner	Abs.	35.3	47.1	42.5	64.3	73.0	18.3	33.8	44.7	14.6	21.2	28.0	11.2	38.7	48.2	21.6
RECOMP-Extr	Ext.	35.8	45.3	2.5	63.5	71.0	2.2	27.6	36.7	2.9	13.8	19.3	3.3	35.2	43.1	2.7
LongLLMLingua	Ext.	32.2	44.0	4.4	66.7	75.2	3.9	34.1	45.3	4.0	28.3	34.8	4.3	40.3	49.8	4.1
Ours (EXIT)	Ext.	36.9	49.4	3.9	67.3	75.9	3.1	37.0	48.3	3.3	28.6	34.5	3.5	42.5	52.0	3.5
Top-20 Documents																
Original Docs	-	39.5	52.5	25.8	69.1	77.6	24.9	38.5	50.0	25.3	28.8	36.8	28.1	44.0	54.2	26
RECOMP-Abst	Abs.	31.5	45.1	4.5	63.4	72.2	3.7	31.3	41.8	4.8	25.4	30.7	4.8	37.9	47.4	4.5
CompAct	Abs.	33.9	45.1	30.8	61.7	70.0	28.1	31.7	40.9	32.0	18.5	23.5	36.5	36.4	44.9	31.9
Refiner	Abs.	32.6	43.5	37.9	62.9	71.4	48.9	31.9	42.2	33.0	22.7	28.7	12.8	37.5	46.5	33.2
RECOMP-Extr	Ext.	33.6	42.7	2.4	63.1	70.3	2.2	25.6	34.5	2.9	12.2	17.4	3.3	33.6	41.2	2.7
LongLLMLingua	Ext.	34.5	46.4	10.9	68.2	76.9	10.4	36.7	48.4	10.6	29.8	36.5	11.0	42.3	52.1	10.7
Ours (EXIT)	Ext.	39.4	52.6	5.1	68.7	77.3	4.3	38.6	50.2	4.7	30.0	36.3	4.8	44.2	54.1	4.7



Table 16: Few-shot QA prompt evaluation measured by EM and F1. Best results for each dataset are highlighted in **bold**, Second best results are highlighted in underline.

Compressor	Type	NQ		TQA		HQA		2WIKI		AVG.	
		EM↑	F1↑	EM↑	F1↑	EM↑	F1↑	EM↑	F1↑	EM↑	F1↑
Llama-3.1-8B-Instruct											
Top-5 Documents											
Original Docs	-	36.9	48.8	61.5	70.3	29.6	39.9	22.2	29.1	37.5	47.0
RECOMP-Abst	Abs.	33.9	45.1	57.8	66.7	27.0	37.2	26.2	32.2	36.2	45.3
CompAct	Abs.	35.0	46.3	60.3	69.2	30.2	40.6	23.9	31.0	37.3	46.8
Refiner	Abs.	34.4	46.0	61.0	70.0	29.6	39.9	23.4	30.0	37.1	46.5
RECOMP-Extr	Ext.	35.9	45.8	58.5	66.1	25.9	35.4	21.2	27.5	35.4	43.7
LongLLMLingua	Ext.	30.7	41.8	60.8	68.8	27.3	37.1	23.0	28.8	35.5	44.1
EXIT (Ours)	Ext.	35.8	47.4	61.0	69.8	29.7	40.3	25.9	32.0	38.1	47.4
Top-20 Documents											
Original Docs	-	39.2	51.4	64.2	73.2	30.6	40.8	24.8	32.4	39.7	49.4
RECOMP-Abst	Abs.	30.2	40.7	59.3	67.6	27.4	37.7	26.8	32.4	35.9	44.6
CompAct	Abs.	35.6	47.0	60.1	69.2	30.7	40.6	21.0	28.3	36.9	46.3
Refiner	Abs.	32.2	43.1	59.6	68.6	27.8	37.9	22.8	29.4	35.6	44.7
RECOMP-Extr	Ext.	34.2	43.7	57.5	64.9	24.6	33.8	19.8	26.0	34.0	42.1
LongLLMLingua	Ext.	33.8	45.2	63.8	72.1	31.0	41.1	25.3	31.6	38.5	47.5
EXIT (Ours)	Ext.	38.8	50.8	63.4	72.3	32.2	43.1	26.6	32.9	40.3	49.8
Llama-3.1-70B-Instruct											
Top-5 Documents											
Original Docs	-	39.5	51.8	68.3	76.5	36.0	46.7	31.4	37.5	43.8	53.1
RECOMP-Abst	Abs.	38.1	50.3	63.4	72.4	30.8	41.2	27.8	33.4	40.0	49.3
CompAct	Abs.	37.9	49.7	67.7	76.0	36.8	47.5	32.2	38.7	43.7	53.0
Refiner	Abs.	38.2	50.2	67.7	76.1	36.0	46.9	30.5	36.6	43.1	52.4
RECOMP-Extr	Ext.	40.1	50.9	68.1	75.6	30.8	40.2	26.5	31.9	41.4	49.7
LongLLMLingua	Ext.	35.2	47.2	69.3	77.2	35.6	46.8	34.7	40.2	43.7	52.8
EXIT (Ours)	Ext.	39.5	51.9	69.6	77.8	38.1	49.4	35.4	41.0	45.7	55.1
Top-20 Documents											
Original Docs	-	42.1	55.0	71.1	79.1	39.4	51.1	35.4	42.4	47.0	56.9
RECOMP-Abst	Abs.	37.2	50.0	65.6	74.0	32.4	43.2	30.3	35.7	41.4	50.8
CompAct	Abs.	37.6	49.1	66.4	74.4	33.6	42.7	25.6	30.5	40.8	49.2
Refiner	Abs.	36.5	47.6	66.6	74.7	33.3	43.3	30.2	35.6	41.6	50.3
RECOMP-Extr	Ext.	38.6	49.1	68.4	75.6	28.9	38.0	24.7	29.9	40.1	48.2
LongLLMLingua	Ext.	37.0	49.1	70.5	78.5	37.9	49.4	35.4	41.1	45.2	54.6
EXIT (Ours)	Ext.	42.5	55.3	71.0	79.0	39.8	51.1	36.5	42.2	47.5	56.9

Table 17: Token distribution analysis across different Top- $k$  scenarios, models, and datasets. Best results for each dataset are highlighted in **bold**, Second best results are highlighted in underline.

Compressor	Type	NQ	TQA	HQA	2WIKI	AVG.
		#token (%) ↓	#token (%) ↓	#token (%) ↓	#token (%) ↓	#token (%) ↓
Top-5 Documents						
Original Docs	-	723.9 (100.0)	730.3 (100.0)	749.2 (100.0)	784.8 (100.0)	734.4 (100.0)
RECOMP-Abst	Abs.	<b>38.0 (5.2)</b>	<b>36.9 (5.0)</b>	<b>63.3 (8.4)</b>	<b>55.1 (7.0)</b>	<b>46.0 (6.2)</b>
CompAct	Abs.	77.5 (10.7)	79.0 (10.8)	77.3 (10.3)	71.4 (9.1)	78 (10.6)
Refiner	Abs.	115.6 (16.0)	103.2 (14.1)	76.6 (10.2)	62.1 (7.9)	98.5 (13.4)
RECOMP-Extr	Ext.	43.9 (6.1)	42.7 (5.8)	90.2 (12.0)	97.9 (12.5)	58.9 (8.0)
LongLLMLingua	Ext.	224.3 (31)	221.9 (30.4)	229.2 (30.6)	234.5 (29.9)	225.1 (30.7)
Ours (EXIT)	Ext.	283.8 (39.2)	211.3 (28.9)	190.3 (25.4)	154.9 (19.7)	228.4 (31.2)
Top-20 Documents						
Original Docs	-	2897.4 (100.0)	2925.2 (100.0)	2996.6 (100.0)	3139.7 (100.0)	2939.8 (100.0)
RECOMP-Abst	Abs.	<b>26.1 (0.9)</b>	<b>38.6 (1.3)</b>	<b>64.5 (2.2)</b>	<b>51.1 (1.6)</b>	<b>43.1 (1.5)</b>
CompAct	Abs.	105.4 (3.6)	102.5 (3.5)	111.8 (3.7)	109.5 (3.5)	106.5 (3.6)
Refiner	Abs.	232.4 (8.0)	176.0 (6.0)	117.3 (3.9)	92.7 (3.0)	175.2 (6.0)
RECOMP-Extr	Ext.	43.4 (1.5)	40.6 (1.4)	89.3 (3.0)	95.7 (3.0)	57.8 (2.0)
LongLLMLingua	Ext.	550.9 (19.0)	553.9 (18.9)	563.3 (18.8)	595.5 (19.0)	556 (18.9)
Ours (EXIT)	Ext.	1001.2 (34.6)	635.0 (21.7)	465.0 (15.5)	367.1 (11.7)	700.4 (23.9)

Table 18: Latency analysis across different Top- $k$  scenarios, models and datasets. Each entry shows compression/reading/total time in seconds. 8B reader experiments were conducted on a single A100-80GB GPU, while 70B reader experiments utilized 4 A100-80GB GPUs in parallel. Best results for each dataset are highlighted in **bold**, Second best results are highlighted in underline.

Compressor	Type	NQ			TQA			HQA			2WIKI			AVG.		
		Comp. ↓	Read ↓	Total ↓	Comp. ↓	Read ↓	Total ↓	Comp. ↓	Read ↓	Total ↓	Comp. ↓	Read ↓	Total ↓	Comp. ↓	Read ↓	Total ↓
Llama-3.1-8B-Instruct																
Top-5 Documents																
Original Docs	-	-	1.03	1.03	-	0.93	0.93	-	0.96	0.96	-	1.12	1.12	-	1.03	1.03
RECOMP-Abst	Abs.	1.13	<b>0.43</b>	1.55	1.06	<b>0.29</b>	1.35	1.92	<b>0.32</b>	2.24	1.73	<b>0.38</b>	2.11	1.46	<b>0.36</b>	1.81
CompAct	Abs.	8.04	<b>0.43</b>	8.47	8.47	0.36	8.83	7.80	0.47	8.26	7.65	<u>0.49</u>	8.14	7.99	0.44	8.43
Refiner	Abs.	27.40	0.70	28.10	10.50	0.40	10.90	6.50	0.40	6.90	5.80	0.50	6.40	12.52	0.55	13.07
RECOMP-Extr	Ext.	<b>0.04</b>	<u>0.50</u>	<b>0.54</b>	<b>0.04</b>	<u>0.31</u>	<b>0.35</b>	<b>0.04</b>	<u>0.37</u>	<b>0.41</b>	<b>0.04</b>	0.59	<b>0.63</b>	<b>0.04</b>	0.44	<b>0.48</b>
LongLLMLingua	Ext.	0.38	0.47	0.85	0.37	0.42	0.79	0.39	0.46	0.84	0.40	0.53	0.93	0.39	0.47	0.86
Ours (EXIT)	Ext.	<u>0.33</u>	<b>0.43</b>	<u>0.76</u>	<u>0.35</u>	0.36	<u>0.71</u>	<u>0.38</u>	0.40	<u>0.78</u>	<u>0.39</u>	0.53	<u>0.92</u>	<u>0.36</u>	<u>0.43</u>	<u>0.79</u>
Llama-3.1-70B-Instruct																
Top-5 Documents																
Original Docs	-	-	3.41	3.41	-	2.85	2.85	-	2.94	2.94	-	3.22	3.22	-	3.11	3.11
RECOMP-Abst	Abs.	<u>1.07</u>	0.63	<u>1.70</u>	1.55	<u>0.31</u>	1.86	2.09	<b>0.35</b>	2.44	2.16	<b>0.44</b>	2.60	1.72	<b>0.43</b>	2.15
CompAct	Abs.	25.58	<u>0.49</u>	26.06	24.08	0.39	24.47	27.02	0.52	27.53	31.63	0.57	32.19	27.08	0.49	27.57
Refiner	Abs.	28.00	0.70	28.70	43.70	0.50	44.20	29.00	0.60	29.60	9.80	1.00	10.80	27.63	0.69	28.32
RECOMP-Extr	Ext.	<b>0.11</b>	0.51	<b>0.62</b>	<b>0.11</b>	<b>0.28</b>	<b>0.39</b>	<b>0.12</b>	0.42	<b>0.54</b>	<b>0.11</b>	0.56	<b>0.68</b>	<b>0.11</b>	<u>0.44</u>	<b>0.55</b>
LongLLMLingua	Ext.	1.76	1.04	2.80	1.78	0.92	2.70	1.83	0.93	2.76	1.89	1.01	2.90	1.81	0.98	2.79
Ours (EXIT)	Ext.	1.33	<b>0.42</b>	1.76	<u>1.37</u>	0.36	<u>1.74</u>	<u>1.45</u>	0.40	<u>1.85</u>	<u>1.51</u>	<u>0.53</u>	<u>2.04</u>	<u>1.42</u>	<b>0.43</b>	<u>1.85</u>
Llama-3.1-70B-Instruct																
Top-5 Documents																
Original Docs	-	-	8.63	8.63	-	7.70	7.70	-	8.30	8.30	-	9.09	9.09	-	8.43	8.43
RECOMP-Abst	Abs.	1.28	3.20	4.48	1.20	<b>2.14</b>	3.34	2.06	<b>2.37</b>	4.43	1.71	<b>2.54</b>	4.24	1.56	<b>2.56</b>	4.12
CompAct	Abs.	8.77	<u>3.11</u>	11.88	8.97	2.72	11.69	8.28	<u>2.73</u>	11.01	8.36	3.23	11.59	8.59	2.95	11.54
Refiner	Abs.	35.90	6.60	42.50	14.70	3.60	18.30	11.30	3.30	14.60	8.00	3.20	11.20	17.48	4.16	21.64
RECOMP-Extr	Ext.	<b>0.04</b>	<b>2.44</b>	<b>2.48</b>	<b>0.04</b>	<u>2.21</u>	<b>2.25</b>	<b>0.04</b>	2.87	<b>2.91</b>	<b>0.05</b>	3.29	<b>3.33</b>	<b>0.04</b>	<u>2.70</u>	<b>2.74</b>
LongLLMLingua	Ext.	0.50	3.87	4.37	0.50	3.40	3.90	0.51	3.52	4.03	0.54	3.72	4.26	0.51	3.63	4.14
Ours (EXIT)	Ext.	<u>0.44</u>	3.50	<u>3.94</u>	<u>0.44</u>	2.66	<u>3.10</u>	<u>0.42</u>	2.88	<u>3.30</u>	<u>0.50</u>	<u>3.03</u>	<u>3.54</u>	<u>0.45</u>	3.02	<u>3.47</u>
Llama-3.1-70B-Instruct																
Top-20 Documents																
Original Docs	-	-	25.78	25.78	-	24.85	24.85	-	25.35	25.35	-	28.09	28.09	-	26.02	26.02
RECOMP-Abst	Abs.	<u>1.13</u>	3.36	<u>4.49</u>	<u>1.58</u>	2.11	<u>3.70</u>	2.24	<u>2.59</u>	4.83	2.07	<b>2.75</b>	4.81	1.76	<u>2.70</u>	<u>4.46</u>
CompAct	Abs.	27.60	<u>3.24</u>	30.84	25.41	2.74	28.15	28.90	3.08	31.99	33.54	<u>2.92</u>	36.45	28.86	2.99	31.86
Refiner	Abs.	32.50	5.40	37.90	47.40	<b>1.50</b>	48.90	31.40	<b>1.50</b>	33.00	9.70	3.10	12.80	30.25	2.90	33.15
RECOMP-Extr	Ext.	<b>0.11</b>	<b>2.27</b>	<b>2.38</b>	<b>0.12</b>	<u>2.10</u>	<b>2.21</b>	<b>0.12</b>	2.77	<b>2.89</b>	<b>0.13</b>	3.18	<b>3.31</b>	<b>0.12</b>	<b>2.58</b>	<b>2.70</b>
LongLLMLingua	Ext.	2.35	8.51	10.86	2.35	8.04	10.38	2.40	8.24	10.65	2.50	8.47	10.97	2.40	8.32	10.71
Ours (EXIT)	Ext.	1.62	3.50	5.12	1.64	2.67	4.31	<u>1.78</u>	2.88	<u>4.65</u>	<u>1.80</u>	2.96	<u>4.75</u>	<u>1.71</u>	3.00	4.71