

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 ¹.

1 Introduction

Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Khandelwal et al., 2020) is the task of enhancing Large Language Models (LLMs) responses 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., 2023b) and improving factual reliability (Jeong et al., 2024; Xia et al., 2024b).

¹We will make our code publicly available upon acceptance of this paper.

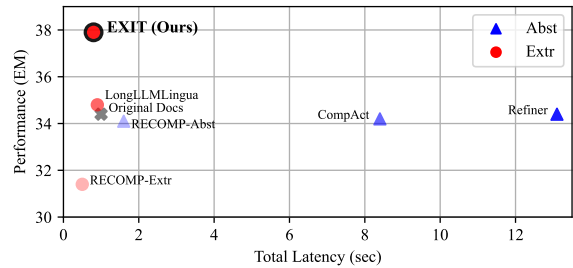


Figure 1: Average QA accuracy (EM) and efficiency (Total Latency) for various compression methods using Contriever-MSMARCO as the retriever and Llama-3.1-8b-Instruct as the reader.

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., 2023). Additionally, irrelevant information in retrieved documents can act as distractors, significantly degrading the overall QA performance (Shi et al., 2023a; Li et al., 2023a; 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 length². Moreover, the context window limitations inherent in LLM architectures set strict upper bounds on the maximum input size.

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

²<https://openai.com/api/pricing/>

069 still have notable drawbacks. Specifically, abstrac- 121
070 tive compression methods—often implemented via 122
071 autoregressive generation—summarize or rewrite 123
072 documents into a single condensed passage (Li 124
073 et al., 2023c; Xu et al., 2024; Yoon et al., 2024; 125
074 Li et al., 2024; Wang et al., 2023), significantly 126
075 increasing end-to-end latency due to their token- 127
076 by-token generation process. For instance, as illus- 128
077 trated in Figure 1, CompAct (Yoon et al., 2024), a 129
078 representative abstractive approach, takes over 8 130
079 seconds to process just five documents for a single 131
080 query, whereas using the original document 132
081 without compression takes only 1 second. 133

082 On the other hand, extractive compression ap- 134
083 proaches can offer a more efficient alternative (Xu 135
084 et al., 2024; Jiang et al., 2023, 2024), by selecting 136
085 relevant textual segments (e.g., sentences, or even 137
086 token-level excerpts) directly from retrieved doc- 138
087 uments. This strategy reduces both compression 139
088 time and overall latency. However, current extrac- 140
089 tive methods have yet to reach their full potential in 141
090 terms of effectiveness (Choi et al., 2021; Pan et al., 142
091 2024). They often rely on rigid selection criteria 143
092 that do not adapt to variations in query complex- 144
093 ity or the quality of retrieved documents, and they 145
094 frequently neglect to fully leverage the broader con- 146
095 text when choosing which tokens or sentences to 147
096 retain. Specifically, as illustrated in Figure 1, while 148
097 extractive approaches such as RECOMP-Extr (Xu 149
098 et al., 2024) achieve minimal compression time, 150
099 their inability to dynamically adjust selection pro- 151
100 cesses results in suboptimal QA performance.

101 Therefore, in this work, we propose a novel com- 152
102 pression framework for RAG, **EX**tract**I**ve **Co**ntex**T** 153
103 compression (EXIT), designed to enhance both ef- 154
104 fectiveness and efficiency by improving efficiency 155
105 through an extractive compression strategy and en- 156
106 hancing effectiveness through dynamic, context- 157
107 aware sentence selection. Specifically, as shown in 158
108 Figure 2, EXIT operates in three stages: (1) split- 159
109 ting retrieved documents into sentences, (2) per- 160
110 forming parallelizable binary classification (“Yes” 161
111 or “No”) on each sentence to assess its relevance 162
112 while considering its full document context, rather 163
113 than evaluating them independently, and (3) recom- 164
114 bining selected sentences while preserving their 165
115 original order. Therefore, as shown in Figure 1, 166
116 EXIT frames context compression as a sentence 167
117 classification problem, enabling it to outperform 168
118 both compression methods and the uncompressed 169
119 baseline in terms of speed. Specifically, it reduces 170
120 processing time from several seconds to about 1

second. Moreover, by leveraging context-aware and 121
adaptive sentence selection, EXIT also surpasses 122
other extractive methods in accuracy. Also, we note 123
that EXIT operates as a plug-and-play module that 124
can be seamlessly integrated into any existing RAG 125
pipeline without architectural modifications. 126

We evaluate EXIT on both single-hop QA tasks 127
(NQ (Kwiatkowski et al., 2019), TQA (Joshi et al., 128
2017)) and multi-hop QA tasks (HQA (Yang et al., 129
2018), 2WikiMultiHopQA (Ho et al., 2020)). Ex- 130
perimental results demonstrate that EXIT not only 131
improves effectiveness over both abstractive and 132
extractive compression baselines but also signif- 133
icantly reduces latency compared to abstractive 134
methods and the uncompressed baseline. 135

Our contributions are as follows: 136

- We identify and address the key weaknesses 137
of existing context compression methods: ab- 138
stractive approaches incur prohibitive latency, 139
while traditional extractive methods rely on 140
rigid, non-adaptive content selection. 141
- We propose **EXIT** (**EX**tract**I**ve **Co**ntex**T** 142
Compression), an extractive compression frame- 143
work that dynamically adjusts to query com- 144
plexity and retrieval quality. 145
- We demonstrate, through extensive experi- 146
ments, that EXIT surpasses previous compres- 147
sion methods and uncompressed retrievals, im- 148
proving QA performance while significantly 149
reducing both token counts and end-to-end 150
latency. 151

2 Related Work 152

Retrieval-Augmented Generation. The RAG 153
pipeline typically follows a naive retrieve-then- 154
generate process, where a single-step retrieval pre- 155
cedes generation (Lewis et al., 2020; Shi et al., 156
2023b; Ram et al., 2023). However, a simple single- 157
step retrieval often fails to rank relevant docu- 158
ments at the top and struggles to handle multi- 159
hop queries requiring multiple pieces of informa- 160
tion. To address these, RAG pipelines have evolved 161
into iterative, recursive, and multi-hop retrieval ap- 162
proaches (Shao et al., 2023; Trivedi et al., 2023; 163
Khattab et al., 2022), which require multiple re- 164
trievals for a single query. While these methods im- 165
prove information coverage, they also increase end- 166
to-end latency from retrieval to generation, reduc- 167
ing the overall efficiency of the pipeline. Moreover, 168

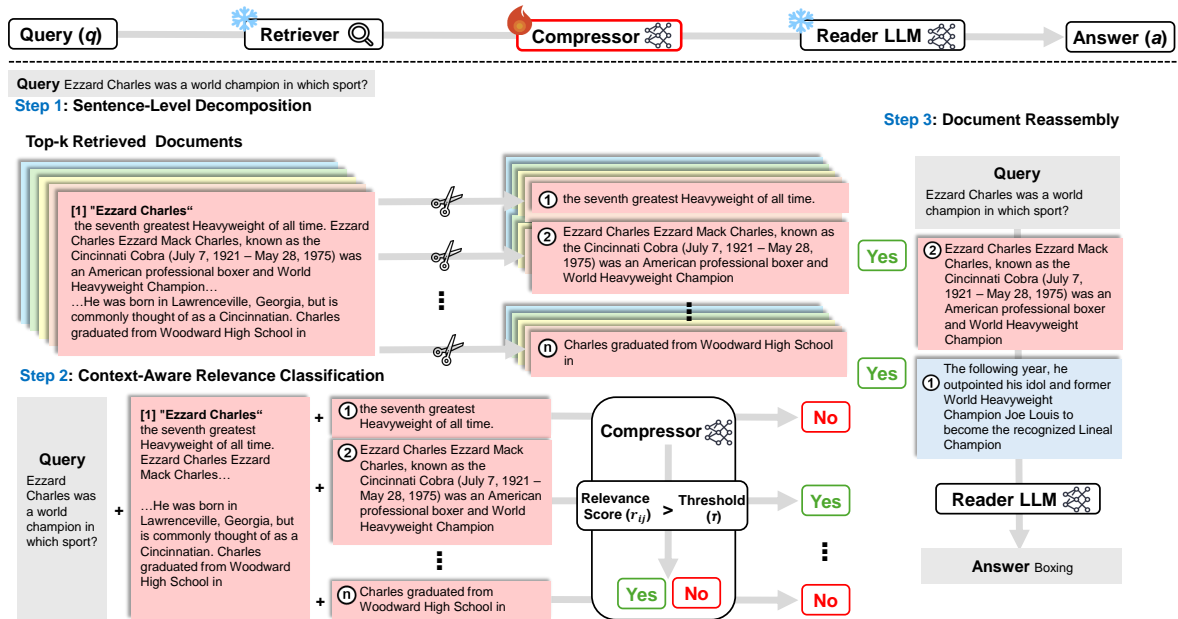


Figure 2: Overview of our framework. First, the retrieved document is split into sentences. Next, each sentence is classified as either “Yes” or “No” using the Compressor. Finally, sentences with scores above the threshold are recombined in their original order to complete the compression.

169 increasing the document length or the number of
170 retrieved documents as alternative solutions to ensure
171 coverage further exacerbates efficiency issues
172 in these complex retrieval approaches. This not
173 only heightens inference costs (Xia et al., 2024a)
174 but also makes it harder to focus on critical details
175 within the documents (Liu et al., 2023). In
176 response, reducing the tokens in the retrieved
177 documents while keeping key information in them
178 has gotten attention to complement the efficiency
179 issue of the RAG pipeline.

180 **Context Compression.** Context compression has
181 emerged as a practical remedy for handling increas-
182 ingly large prompt lengths in RAG pipelines. Ex-
183 isting approaches commonly fall into soft or hard
184 compression. Soft compression focuses on shorten-
185 ing embedding vectors at the token level (Wingate
186 et al., 2022; Mu et al., 2023; Ge et al., 2024; Cheva-
187 lier et al., 2023; Cheng et al., 2024), but requires
188 extensive training and architectural changes, mak-
189 ing it unsuitable for black-box LLMs.

190 Hard compression, by contrast, removes non-
191 essential textual content directly (Li et al., 2023d;
192 Jiang et al., 2023), offering a plug-and-play so-
193 lution compatible even with API-based models
194 such as ChatGPT (OpenAI, 2023). Hard compres-
195 sion techniques are further divided into abstrac-
196 tive and extractive methods. Abstractive methods
197 employ autoregressive models to generate query-
198 focused summaries, thus drastically reducing to-
199 ken counts at the cost of additional latency and

200 potential hallucinations (Zhao et al., 2020). For
201 example, RECOMP-Abst (Xu et al., 2024) uses a
202 T5-based summarizer for token reduction but re-
203 quires dataset-specific training and slows inference.
204 CompAct (Yoon et al., 2024) and Refiner (Li et al.,
205 2024) take this approach further by leveraging even
206 larger LLMs with 7B parameters, compounding la-
207 tency issues and increasing resource demands.

208 Extractive methods select salient segments (e.g.,
209 sentences, or tokens) directly from the retrieved
210 documents. This avoids the autoregressive bottle-
211 neck and mitigates hallucinations. RECOMP-Extr
212 (Xu et al., 2024) is one such example, but its static
213 and context-agnostic selection of only a few sen-
214 tences per document limits its performance. Simi-
215 larly, token-level approaches such as LLMLingua
216 family (Li et al., 2023d; Jiang et al., 2023; Pan
217 et al., 2024; Jiang et al., 2024) can distort seman-
218 tic coherence by removing key entities or splitting
219 essential facts.

220 In short, abstractive methods offer strong com-
221 pression but suffer from latency and potential hal-
222 lucinations, while extractive methods are often
223 rigid and lack context awareness. Our work ad-
224 dresses these limitations by proposing a paralleliz-
225 able, context-aware extractive compression frame-
226 work. It adaptively selects sentences at scale, pre-
227 serving semantic integrity and efficiently balancing
228 accuracy with speed, even in complex, multi-step
229 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. For effective and efficient compression in the RAG pipeline, three key criteria should be satisfied: (1) D' should contain fewer tokens, as fewer tokens lead to shorter answer generation times (i.e., reading time); (2) D' should retain the essential evidence required to answer the query, ensuring effectiveness; and (3) the compression process should be sufficiently fast enough to avoid significantly increasing the overall end-to-end inference latency.

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 $d_i \in D$, we produce 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 effectively retained.

Context-Aware Relevance Classification. To effectively and efficiently filter sentences in D that contain key evidence for answering the question, we design two key components for sentence relevance evaluation: **document consideration** and **single-token prediction**. First, incorporating the entire document d_i is essential, as key evidence may be distributed throughout the context rather than confined to a single sentence, ensuring no critical information is missed and enabling effective context compression. Furthermore, as multiple-token generation proposed in the previous work (Yoon et al., 2024; Li et al., 2024) could compromise the overall efficiency of the RAG pipeline, we design the lightweight relevance calculation with only single-token prediction with “Yes” and “No” from the query q , document d_i , and sentence s_{ij} . In detail, for each candidate sentence s_{ij} , the evaluation model calculates the relevance score r_{ij} of the sentence s_{ij} for a given query q and document d_i as follows:

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

where $P(\cdot|\cdot)$ denotes the likelihood of each token from the evaluation model. This relevance calculation is parallelized across multiple sentences, allowing them to be evaluated simultaneously.

Then, among the sentences in D , EXIT selects sentences with a relevance score exceeding a pre-defined threshold τ , as high relevance scores indicate that a sentence contains critical information. Notably, this selection process results in **an adaptive number of sentences** in the compressed set D' , rather than a fixed amount. This adaptive approach aligns with prior work (Jeong et al., 2024), recognizing that the complexity of queries and the amount of key information vary across queries. As a result, our sentence selection strategy enables effective compression while ensuring all key evidence is included in the compressed set.

Document Reassembly. As shown in Step 3 of Figure 2, EXIT reconstructs the compressed document D' by concatenating only the selected sentences in their original order. Following Hwang et al. (2024), preserving the canonical sentence

sequence maintains logical flow and contextual coherence. This approach ensures that the resulting compressed document remains understandable and supports 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) content that directly provides the needed evidence, (2) passages that appear relevant but lack crucial details, and (3) texts that are entirely irrelevant.

Data Sampling. From the annotated dataset, we draw positive examples from sentences explicitly marked as necessary for producing the correct answer. Within the same documents, we select hard-negative examples—sentences that appear related to the query but do not contain the required evidence—thereby simulating plausible but incomplete retrieval scenarios. Additionally, we sample random negatives from unrelated queries, ensuring the classifier learns to dismiss off-topic content. By maintaining a balanced mix of positives, hard negatives, and random negatives, we create a training set that captures a wide spectrum of retrieval conditions.

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 A.

Datasets. We evaluate on both single-hop and multi-hop question answering datasets: **Natu-**

ralQuestions (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 following context compression approaches: **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 evaluate our approach using 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 is defined as the total time encompassing both the compression and generation steps, as the retrieval step remains consistent across all methods.

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	-	<u>34.6</u>	<u>47.1</u>	1.0	58.8	68.6	0.9	28.1	38.6	1.0	16.1	24.9	1.1	34.4	<u>44.8</u>	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	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	<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.	35.9	47.8	<u>0.8</u>	60.8	69.9	<u>0.7</u>	30.6	41.5	<u>0.8</u>	24.2	30.8	<u>0.9</u>	37.9	47.5	<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	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	<u>66.7</u>	<u>75.2</u>	3.9	<u>34.1</u>	<u>45.3</u>	4.0	<u>28.3</u>	34.8	4.3	<u>40.3</u>	49.8	4.1
EXIT (Ours)	Ext.	36.9	49.4	<u>3.9</u>	67.3	75.9	<u>3.1</u>	37.0	48.3	<u>3.3</u>	28.6	<u>34.5</u>	<u>3.5</u>	42.5	52.0	<u>3.5</u>

5 Main Results

Table 1 summarizes our evaluation results across multiple datasets and compression strategies. 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 approaches.

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 uncompressed 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, make 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 B.

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 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.

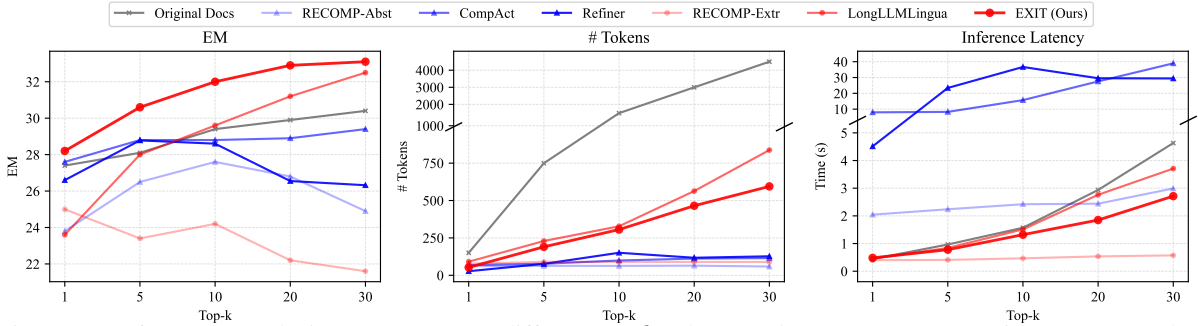


Figure 3: Performance analysis on HQA across different top- k values (1, 5, 10, 20, 30), comparing accuracy, token retention, and inference latency between baselines and our method. All experiments were conducted on a single A100-80GB GPU.

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

Figure 4: Confusion matrices (row-normalized) for context-aware relevance classification on HQA (in-domain) and 2WIKI (out-of-domain).

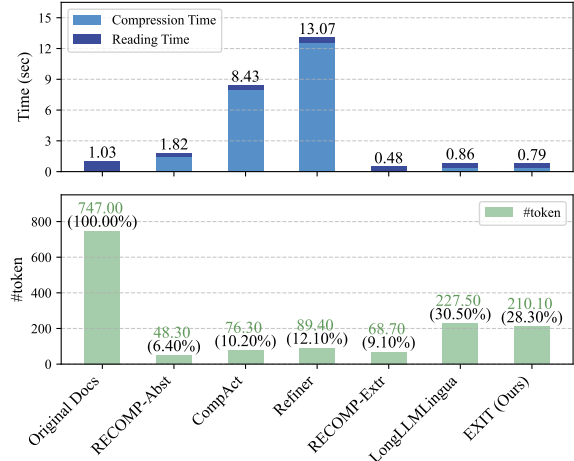


Figure 5: Comparison of compression and reading latency across baselines and our method in QA setting. Experiments were conducted on a single A100 GPU.

6.2 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 4. 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. 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.

6.3 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 read time and compression time, along with an analysis of the average number of tokens in compressed documents, as shown in Figure 5. 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 compression 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 method, EXIT, aligns closely with these objectives, offering a practical approach to context compression for RAG.

6.4 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 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)	31.6	42.6	195.1
Pos + H-Neg	30.0	41.3	286.8
Pos + Neg	29.8	40.9	404.6
w/o Adaptive Sentence Selection	29.4	40.7	91.0
w/o Passage as context	30.4	42.3	157.4

533 bins positive, hard negative, and random negative
 534 samples to mirror the diversity of real-world re-
 535 trieval scenarios. Compared to using only subsets
 536 of these sample types, the comprehensive strategy
 537 improves EM by 1.6 points over using only hard
 538 negatives and by 1.8 points over only random nega-
 539 tives. Relying solely on positive and random nega-
 540 tives led to excessive token retention, while depend-
 541 ing only on hard negatives diminished the model’s
 542 ability to filter out spurious retrieval noise.

543 **Adaptive Sentence Selection.** We compare our
 544 adaptive selection mechanism to fixed-length se-
 545 lection. Although fixed-length selection achieves
 546 the lowest token count (91.0), it reduces EM by 2.2
 547 points and F1 by 1.9 points. This underscores the
 548 importance of adaptively selecting sentences based
 549 on the complexity of the query and the retrieved
 550 documents, rather than using a static cap.

551 **Context-Aware Extraction.** We assess the impact
 552 of incorporating full document context when eval-
 553 uating each sentence’s relevance. Removing sur-
 554 rounding context saves 38 tokens but lowers EM
 555 by 1.2 points, indicating that broader contextual
 556 awareness is crucial for maintaining answer accu-
 557 racy, even if it slightly increases token count.

558 These findings confirm that a balanced training
 559 data strategy enhances robustness, that adaptive
 560 sentence selection ensures efficiency, and that full
 561 document consideration preserves accuracy. The
 562 classification performances under each ablation set-
 563 ting are reported in Appendix B.5.

564 6.5 Impact of Compressor Model Size and 565 Compression Strategy

566 **Model Size Considerations.** Figure 6 presents an
 567 ablation study examining how different base mod-
 568 els influence EM scores and total latency. Note
 569 that all models trained within EXIT achieve supe-
 570 rior accuracy compared to uncompressed baselines
 571 and fast compression under 2 seconds. Specifically,
 572 Gemma-2B, our base classifier model, achieves a

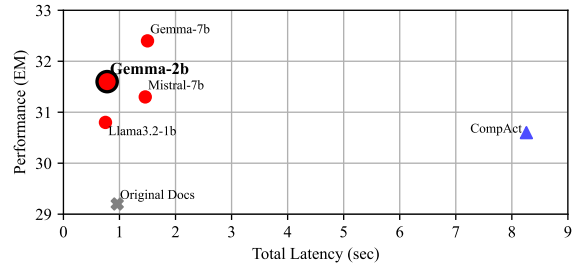


Figure 6: Ablation on HQA comparing EM scores and latency for different model configurations within EXIT (red dot). CompAct and Original Docs are included as indicators. Experiments on a single A100-80GB GPU.

573 favorable balance between effectiveness and effi-
 574 ciency, delivering 31.6 EM points in 0.78s. When
 575 moving to the largest model, Gemma-7B, the high-
 576 est accuracy (32.4 EM) is achieved, yet it also in-
 577 flates latency to 1.50s, slightly exceeding the time
 578 of uncompressed documents. These results sug-
 579 gest that scaling up model parameters can improve
 580 performance but may also compromise latency ben-
 581 efits, emphasizing the flexibility of our proposed
 582 framework by selecting an appropriate model de-
 583 pending on the user requirement.

584 **Abstractive vs. Extractive Compression.** Fig-
 585 ure 6 also compares our extractive approach, EXIT,
 586 against CompAct, a 7B-scale abstractive compres-
 587 sor. Using Mistral-7B as the base model for both
 588 methods, EXIT (31.3 EM, 1.46s) significantly out-
 589 performs CompAct (30.6 EM, 8.26s) in terms of
 590 latency and maintains competitive accuracy. This
 591 stark difference underscores that the compression
 592 strategy, not the just model size, heavily influences
 593 efficiency. By relying on extraction rather than it-
 594 erative summarization, EXIT capitalizes on a large-
 595 scale model to preserve high accuracy without in-
 596 ccurring prohibitively long inference time.

597 7 Conclusion

598 We present EXIT, an efficient context compression
 599 framework for RAG systems that leverages parallel
 600 processing and context-aware, adaptive sentence
 601 selection. Our experiments demonstrate that EXIT
 602 achieves superior performance across both single-
 603 hop and multi-hop QA tasks while maintaining
 604 practical inference speeds. Despite being trained
 605 only on HQA, EXIT shows a strong zero-shot gen-
 606 eralization ability and proves effective across a
 607 wide range of open-source models of varying sizes.
 608 These results suggest that efficient parallel extrac-
 609 tion with smaller models can outperform larger ab-
 610 stractive approaches, offering a practical solution
 611 for real-world RAG applications.

612 Limitation

613 Our current approach relies on explicit sentence-
614 level annotations to train the classifier. While these
615 annotations were obtained manually in our experi-
616 ments, they could potentially be automated through
617 alternative means, such as by GPT-4 supervision or
618 signals derived from the reader itself. We have not
619 yet explored these automated annotation strategies,
620 but doing so remains a promising avenue for future
621 work. Additionally, our study primarily focuses on
622 a general-domain setting, leaving questions about
623 the classifier’s performance in specialized domains
624 unanswered. Investigating how well our approach
625 generalizes to domain-specific or highly special-
626 ized corpora presents another valuable direction for
627 future research. Lastly, we focus on a single-step
628 RAG pipeline, where retrieval occurs only once, ex-
629 cluding more complex pipelines (Shao et al., 2023;
630 Trivedi et al., 2023; Khattab et al., 2022). However,
631 our proposed framework, EXIT, is orthogonal to
632 these approaches and can be seamlessly integrated
633 by compressing the retrieved documents from each
634 retrieval step.

635 Ethics Statement

636 This work enhances RAG-based QA without gener-
637 ating new content beyond what is retrieved. How-
638 ever, biases and inaccuracies in the source docu-
639 ments can still propagate through our compres-
640 sion process. Ensuring the reliability, fairness, and
641 proper curation of underlying corpora is essential
642 for ethical deployment. Future efforts should in-
643 tegrate bias detection, provenance tracking, and
644 user-centric evaluations to promote more transpar-
645 ent and equitable real-world applications.

646 Acknowledgments

647 References

648 Xin Cheng, Xun Wang, Xingxing Zhang, Tao Ge, Si-
649 Qing Chen, Furu Wei, Huishuai Zhang, and Dongyan
650 Zhao. 2024. xrag: Extreme context compression
651 for retrieval-augmented generation with one token.
652 *CoRR*, abs/2405.13792.

653 Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and
654 Danqi Chen. 2023. *Adapting language models to*
655 *compress contexts*. In *Proceedings of the 2023 Con-*
656 *ference on Empirical Methods in Natural Language*
657 *Processing, EMNLP 2023, Singapore, December 6-*
658 *10, 2023*, pages 3829–3846. Association for Compu-
659 tational Linguistics.

660 Eunsol Choi, Jennimaria Palomaki, Matthew Lamm,
661 Tom Kwiatkowski, Dipanjan Das, and Michael

Collins. 2021. *Decontextualization: Making sen-*
662 *tences stand-alone*. *Trans. Assoc. Comput. Linguis-*
663 *tics*, 9:447–461. 664

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, 665
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, 666
Akhil Mathur, Alan Schelten, Amy Yang, Angela 667
Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, 668
Archi Mitra, Archie Sravankumar, Artem Korenev, 669
Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien 670
Rodriguez, Austen Gregerson, Ava Spataru, Bap- 671
tiste Rozière, Bethany Biron, Binh Tang, Bobbie 672
Chern, Charlotte Caucheteux, Chaya Nayak, Chloe 673
Bi, Chris Marra, Chris McConnell, Christian Keller, 674
Christophe Touret, Chunyang Wu, Corinne Wong, 675
Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al- 676
lonsius, Daniel Song, Danielle Pintz, Danny Livshits, 677
David Esiobu, Dhruv Choudhary, Dhruv Mahajan, 678
Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, 679
Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, 680
Emily Dinan, Eric Michael Smith, Filip Radenovic, 681
Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Geor- 682
gia Lewis Anderson, Graeme Nail, Grégoire Mialon, 683
Guan Pang, Guillem Cucurell, Hailey Nguyen, Han- 684
nah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, 685
Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan 686
Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan 687
Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, 688
Jeet Shah, Jelmer van der Linde, Jennifer Billock, 689
Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, 690
Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, 691
Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph 692
Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, 693
Kalyan Vasuden Alwala, Kartikeya Upasani, Kate 694
Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and 695
et al. 2024. *The llama 3 herd of models*. *CoRR*, 696
abs/2407.21783. 697

Tao Ge, Jing Hu, Lei Wang, Xun Wang, Si-Qing Chen, 698
and Furu Wei. 2024. *In-context autoencoder for con-*
699 *text compression in a large language model*. In *The*
700 *Twelfth International Conference on Learning Rep-*
701 *resentations, ICLR 2024, Vienna, Austria, May 7-11,*
702 *2024*. OpenReview.net. 703

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, 704
and Akiko Aizawa. 2020. *Constructing A multi-hop*
705 *QA dataset for comprehensive evaluation of reason-*
706 *ing steps*. In *Proceedings of the 28th International*
707 *Conference on Computational Linguistics, COLING*
708 *2020, Barcelona, Spain (Online), December 8-13,*
709 *2020*, pages 6609–6625. International Committee on
710 Computational Linguistics. 711

Taeho Hwang, Soyeong Jeong, Sukmin Cho, SeungY- 712
oon Han, and Jong Park. 2024. *DSLRL: Document*
713 *refinement with sentence-level re-ranking and recon-*
714 *struction to enhance retrieval-augmented generation*.
715 In *Proceedings of the 3rd Workshop on Knowledge*
716 *Augmented Methods for NLP*, pages 73–92, Bangkok,
717 Thailand. Association for Computational Linguistics. 718

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebas- 719
tian Riedel, Piotr Bojanowski, Armand Joulin, and 720

721	Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning . <i>Trans. Mach. Learn. Res.</i> , 2022.	<i>Symposium on Operating Systems Principles, SOSP 2023, Koblenz, Germany, October 23-26, 2023</i> , pages 611–626. ACM.	779
722			780
723			781
724	Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju Hwang, and Jong C. Park. 2024. Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity . <i>arXiv:2403.14403</i> , abs/2403.14403.	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. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks . 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> .	782
725			783
726			784
727			785
728			786
729	Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023. Llmlingua: Compressing prompts for accelerated inference of large language models . 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.		787
730			788
731			789
732			790
733		Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix X. Yu, and Sanjiv Kumar. 2023a. Large language models with controllable working memory . 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.	791
734			792
735			793
736			794
737	Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2024. Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression . 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 1658–1677. Association for Computational Linguistics.		795
738			796
739			797
740			798
741		Junyi Li, Xiaoxue Cheng, Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023b. Halueval: A large-scale hallucination evaluation benchmark for large language models . 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.	799
742			800
743			801
744			802
745			803
746	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension . 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.		804
747			805
748			806
749			807
750			808
751			809
752			810
753			811
754			812
755	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through memorization: Nearest neighbor language models . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.		813
756			814
757			815
758			816
759			817
760	Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. 2022. Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive NLP . <i>CoRR</i> , abs/2212.14024.		818
761			819
762			820
763			821
764			822
765	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris Alberti, 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. Natural questions: a benchmark for question answering research . <i>Trans. Assoc. Comput. Linguistics</i> , 7:452–466.		823
766			824
767			825
768			826
769			827
770			828
771			829
772			830
773			831
774	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention . In <i>Proceedings of the 29th</i>		832
775			833
776			834
777			835
778			

836	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. Gemma: Open models based on gemini research and technology . <i>CoRR</i> , abs/2403.08295.	
847	Jesse Mu, Xiang Li, and Noah D. Goodman. 2023. Learning to compress prompts with gist tokens . 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> .	
853	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A human generated machine reading comprehension dataset . 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.	
863	OpenAI. 2023. GPT-4 technical report . <i>arXiv:2303.08774</i> , abs/2303.08774.	
865	Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor Rühle, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and Dongmei Zhang. 2024. Llmlingua-2: Data distillation for efficient and faithful task-agnostic prompt compression . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 963–981. Association for Computational Linguistics.	
874	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models . <i>Trans. Assoc. Comput. Linguistics</i> , 11:1316–1331.	
879	Stephen E. Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: BM25 and beyond . <i>Found. Trends Inf. Retr.</i> , 3(4):333–389.	
882	Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 9248–9274. Association for Computational Linguistics.	
890	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H. Chi, Nathanael Schärli, and Denny Zhou. 2023a. Large language models can	
	be easily distracted by irrelevant context . In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 31210–31227. PMLR.	893 894 895 896 897
	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023b. REPLUG: retrieval-augmented black-box language models . <i>arXiv.2301.12652</i> , abs/2301.12652.	898 899 900 901 902
	Harsh Trivedi, Niranjana Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 10014–10037. Association for Computational Linguistics.	903 904 905 906 907 908 909 910 911
	Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md. Rizwan Parvez, and Graham Neubig. 2023. Learning to filter context for retrieval-augmented generation . <i>arXiv.2311.08377</i> , abs/2311.08377.	912 913 914 915
	David Wingate, Mohammad Shoeybi, and Taylor Sorensen. 2022. Prompt compression and contrastive conditioning for controllability and toxicity reduction in language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 5621–5634. Association for Computational Linguistics.	916 917 918 919 920 921 922 923
	Zhenyu Wu, Chao Shen, and Meng Jiang. 2024. Instructing large language models to identify and ignore irrelevant conditions . <i>arXiv.2403.12744</i> , abs/2403.12744.	924 925 926 927
	Heming Xia, Zhe Yang, Qingxiu Dong, Peiyi Wang, Yongqi Li, Tao Ge, Tianyu Liu, Wenjie Li, and Zhifang Sui. 2024a. Unlocking efficiency in large language model inference: A comprehensive survey of speculative decoding . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 7655–7671. Association for Computational Linguistics.	928 929 930 931 932 933 934 935 936
	Peng Xia, Kangyu Zhu, Haoran Li, Hongtu Zhu, Yun Li, Gang Li, Linjun Zhang, and Huaxiu Yao. 2024b. RULE: reliable multimodal RAG for factuality in medical vision language models . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 1081–1093. Association for Computational Linguistics.	937 938 939 940 941 942 943 944
	Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2024. RECOMP: Improving retrieval-augmented LMs with context compression and selective augmentation . In <i>The Twelfth International Conference on Learning Representations</i> .	945 946 947 948 949

950 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-
951 gio, William W. Cohen, Ruslan Salakhutdinov, and
952 Christopher D. Manning. 2018. [Hotpotqa: A dataset](#)
953 [for diverse, explainable multi-hop question answering](#).
954 In *Proceedings of the 2018 Conference on Em-*
955 *pirical Methods in Natural Language Processing,*
956 *Brussels, Belgium, October 31 - November 4, 2018,*
957 pages 2369–2380. Association for Computational
958 Linguistics.

959 Chanwoong Yoon, Taewhoo Lee, Hyeon Hwang, Min-
960 byul Jeong, and Jaewoo Kang. 2024. [Compact: Com-](#)
961 [pressing retrieved documents actively for question](#)
962 [answering](#). *CoRR*, abs/2407.09014.

963 Zheng Zhao, Shay B. Cohen, and Bonnie Webber. 2020.
964 [Reducing quantity hallucinations in abstractive sum-](#)
965 [marization](#). *CoRR*, abs/2009.13312.

Appendix

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

A 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.

A.1 Training Configuration

We trained the compressor model using the following settings:

- Batch size: 8 per device
- Gradient accumulation steps: 8
- Learning rate: $1e-5$
- Weight decay: 0.1
- Warmup ratio: 0.03
- Training epochs: 1

Memory Optimization:

- Optimizer: paged_adamw_8bit
- Quantization: 4-bit with float16
- LoRA configuration: Rank=64, Scaling=32, Dropout=0.05

A.2 Model Selection and Training Time

Model selection was based on validation loss. Training required approximately 90 hours on our cluster.

A.3 Data Processing

Table 3: 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

We used SpaCy to segment documents into sentences. Table 3 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), 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.

A.4 Inference Settings

For inference, we set:

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

A.5 Prompt Templates

Table 4: 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”.

Table 5: A prompt template used by the reader model for the QA task.

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

Full prompt templates for compression and QA tasks are provided in Tables 4 and 5, respectively.

A.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: Performance comparison between in-domain (HQA) and out-of-domain (2WIKI) datasets using BM25 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
<i>Top-5 Documents</i>							
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	63.0 (8.3)	25.0	30.6	55.7 (7.1)
CompAct	Abs.	30.6	41.2	<u>76.9 (10.2)</u>	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	<u>62.9 (8.0)</u>
RECOMP-Extr	Ext.	28.2	37.5	93.1 (12.3)	11.8	19.1	<u>99.1 (12.5)</u>
LongLLMLingua	Ext.	28.6	39.7	230.3 (30.5)	22.2	27.4	236.6 (30.0)
EXIT (Ours)	Ext.	33.4	44.5	177.8 (23.5)	<u>24.4</u>	<u>29.1</u>	138.2 (17.5)
<i>Top-20 Documents</i>							
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	69.7 (2.3)	22.8	28.2	52.3 (1.7)
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	<u>96.7 (3.1)</u>
LongLLMLingua	Ext.	31.0	40.7	558.0 (18.5)	<u>23.6</u>	28.3	593.7 (19.0)
EXIT (Ours)	Ext.	35.2	46.9	411.3 (13.7)	27.2	32.4	312.7 (10.0)

B Additional Experimental Results and Analyses

Table 13 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 14 summarizes the results, where we randomly selected five training examples per dataset as demonstrations. Table 15 provides comprehensive token count statistics, and Table 16 breaks down end-to-end latency.

B.1 Performance with Sparse Retrieval

To assess EXIT’s robustness with different retrieval architectures, we evaluate it with BM25, a sparse retrieval method. Table 6 compares EXIT’s performance on HQA (in-domain) and 2WIKI (out-of-domain) under Top-5 and Top-20 retrieval settings.

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.

Table 7: 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
<i>Top-5 Documents</i>							
Original Docs	-	37.2	48.6	735.3 (100.0)	31.2	35.3	764.5 (100.0)
RECOMP-Abst	Abs.	29.4	40.2	62.8 (8.5)	23.8	27.8	53.5 (7.0)
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.	38.2	50.4	191.2 (26.0)	31.8	35.8	145.6 (19.0)
<i>Top-20 Documents</i>							
Original Docs	-	39.6	51.8	2940.5 (100.0)	40.0	43.8	3066.2 (100.0)
RECOMP-Abst	Abs.	33.6	44.2	62.7 (2.1)	26.6	32.1	48.9 (1.6)
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	86.2 (2.9)	23.2	27.2	91.0 (3.0)
LongLLMLingua	Ext.	38.8	49.4	549.5 (18.7)	35.4	39.6	581.5 (19.0)
EXIT (Ours)	Ext.	39.4	<u>50.1</u>	453.6 (15.4)	35.6	<u>40.3</u>	346.7 (11.3)

B.2 Performance with Proprietary Model

We further examine EXIT’s effectiveness using GPT-4o as the reader. Table 7 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 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.

B.3 Impact of Threshold τ

We analyze EXIT’s sensitivity to the relevance threshold τ . Figure 7 shows EXIT’s performance across various τ values.

EXIT remains stable over a wide threshold range, with strong results between $\tau=0.3-0.5$. At $\tau=0.3$,

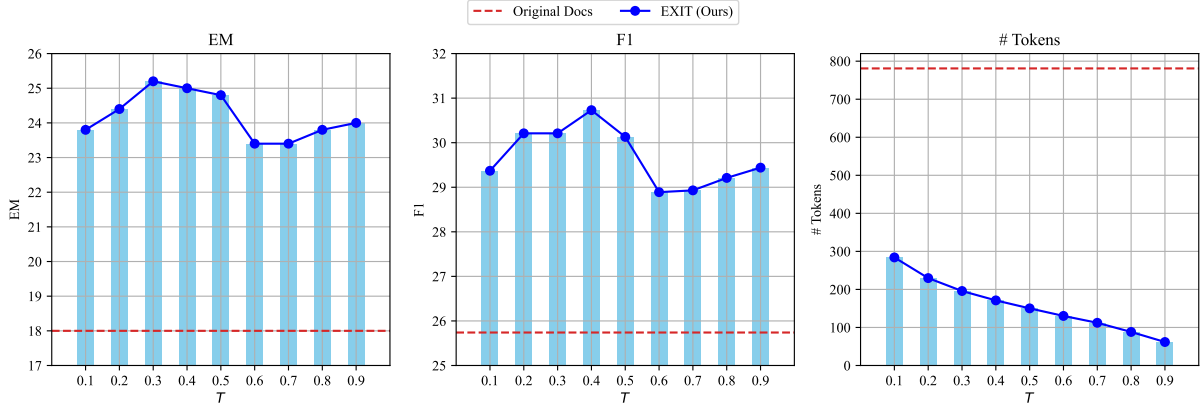


Figure 7: Changes in EM, F1 score, and token count as the threshold T for retaining sentences is adjusted.

Table 8: Classification performance for Yes/No labels, including overall accuracy.

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

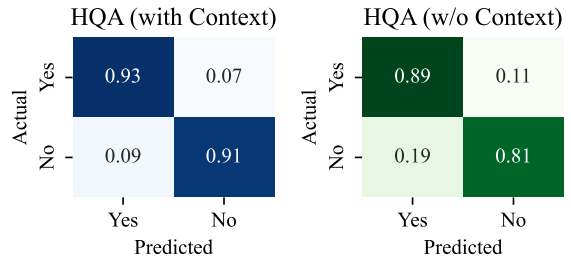


Figure 8: Row-normalized confusion matrices comparing classification performance with (left) and without (right) contextual information on HQA. The availability of context improves the model’s ability to accurately distinguish relevant (“Yes”) from non-relevant (“No”) sentences.

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.

B.4 Analysis of Classification Performance Across Negative Sample Types

Table 8 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, indi-

cating 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.

B.5 Classification Performance under Ablation Setting.

B.5.1 Analysis of Training Data Composition

Figure 9 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

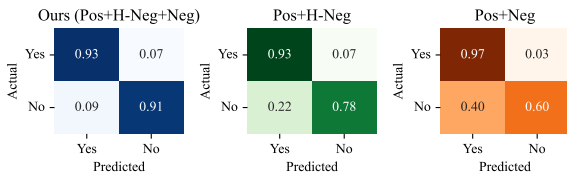


Figure 9: Row-normalized confusion matrices for classification performance under different training data conditions: Ours (Pos+H-Neg+Neg), Pos+H-Neg, and Pos+Neg. Each matrix compares the predicted (“Yes”/“No”) labels against the actual labels.

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.

B.5.2 Impact of Context on Classification Performance

We evaluate the classifier’s performance with and without broader passage-level context. Figure 8 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.

B.6 Training Data Ablation Analysis

Table 9 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.

Table 9: Training data ablation comparison. Best results per metric are highlighted in **bold**.

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

B.7 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 10, 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 11), 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 12 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.

1205 In contrast, Ours avoids hallucinations and retains
1206 the answer’s entity in a concise, coherent form.

1207 LongLLMLingua’s token-level filtering ap-
1208 proach yields unreadable text and removes the
1209 essential “Romania” entity, preventing the reader
1210 from generating the correct answer. Ours, on the
1211 other hand, maintains semantic coherence and in-
1212 cludes the correct entity, allowing the reader to
1213 produce the correct answer without interference.

1214 These case studies highlight the advantages of
1215 Ours: it eliminates distractors, preserves critical
1216 entities, and maintains semantic integrity. Conse-
1217 quently, the reader consistently arrives at correct
1218 answers with fewer tokens and no hallucinations.

Table 10: Case study comparing compressed contexts and answers between Original Docs and Ours.

	Original Docs	Ours
Query	Crème Anglaise is the French version of which English dessert item?	
Context	<p>[1] Creme anglaise Crème anglaise Crème anglaise (French for ""English cream"") is a light pouring Cus-tard used as a dessert cream or sauce. It is a mix of sugar, egg yolks, oil, and hot milk ...</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 Eggnog 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>[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 Cus-tard 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>
Answer	Custard	
Predict	Eggnog	Custard

Table 11: 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
Query	Which 90s sci fi series with James Belushi was based on Bruce Wagner's comic strip of the same name?	
Context	<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, " Probe ", 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>	<i>No context</i> (completely pruned)
Answer	Wild Palms	
Predict	Probe	Wild Palms

Table 12: Combined case studies comparing compressed contexts and answers across different methods.

Case Study 1	CompAct	Ours
Query	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)?	
Context	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).
Answer	Pasadena, Texas	
Predict	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
Query	Who said that there are only three true sports, bullfighting, mountaineering, and motor racing, the rest are merely games?	
Context	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.
Answer	Ernest Hemingway	
Predict	There is no information provided about who said that.	Ernest Hemingway
Case Study 3	LongLLMLingua	Ours
Query	Which was the only Eastern Bloc country to participate in the 1984 LA Olympics?	
Context	984 Summer Olympicscottthe majority of Soviet Bloc will not particip Games,ausescu's Roman expected.984 Summer boyc The boyc 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 initi boycott on May 8, 1984. Bootting countries organized major event, the Friendship Games, in July and August 1984 Although the boycott by the Soviet Union affected number Olympic Sum 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.
Answer	Romania	
Predict	China	Romania

Table 13: 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																
<i>Top-5 Documents</i>																
Original Docs	-	<u>34.6</u>	<u>47.1</u>	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.	<u>34.6</u>	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	<u>59.4</u>	68.0	0.8	28.0	38	<u>0.8</u>	21.5	27.4	<u>0.9</u>	<u>34.8</u>	43.7	0.9
Ours (EXIT)	Ext.	35.9	47.8	<u>0.8</u>	60.8	69.9	<u>0.7</u>	30.6	41.5	<u>0.8</u>	24.2	30.8	<u>0.9</u>	37.9	47.5	<u>0.8</u>
<i>Top-20 Documents</i>																
Original Docs	-	<u>36.6</u>	<u>49.5</u>	3.4	62.0	<u>71.7</u>	2.9	<u>29.9</u>	40.5	2.9	18.8	27.9	3.2	36.8	<u>47.4</u>	3.1
RECOMP-Abst	Abs.	26.9	38.3	<u>1.7</u>	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	<u>62.4</u>	71.2	2.7	31.2	41.4	2.8	<u>24.1</u>	<u>30.1</u>	2.9	<u>37.8</u>	46.9	2.8
Ours (EXIT)	Ext.	38.1	50.8	1.8	62.8	72.0	<u>1.7</u>	32.9	44.0	<u>1.8</u>	25.5	32.3	<u>2.0</u>	39.8	49.8	<u>1.8</u>
Llama-3.1-70B-Instruct																
<i>Top-5 Documents</i>																
Original Docs	-	35.6	<u>48.0</u>	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.	<u>35.8</u>	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	<u>66.7</u>	<u>75.2</u>	3.9	<u>34.1</u>	45.3	4.0	<u>28.3</u>	34.8	4.3	<u>40.3</u>	<u>49.8</u>	4.1
Ours (EXIT)	Ext.	36.9	49.4	<u>3.9</u>	67.3	75.9	<u>3.1</u>	37.0	48.3	<u>3.3</u>	28.6	<u>34.5</u>	<u>3.5</u>	42.5	52.0	<u>3.5</u>
<i>Top-20 Documents</i>																
Original Docs	-	39.5	<u>52.5</u>	25.8	69.1	77.6	24.9	<u>38.5</u>	<u>50.0</u>	25.3	28.8	36.8	28.1	<u>44.0</u>	54.2	26
RECOMP-Abst	Abs.	31.5	45.1	<u>4.5</u>	63.4	72.2	<u>3.7</u>	31.3	41.8	4.8	25.4	30.7	<u>4.8</u>	37.9	47.4	<u>4.5</u>
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	<u>29.8</u>	36.5	11.0	42.3	52.1	10.7
Ours (EXIT)	Ext.	<u>39.4</u>	52.6	5.1	<u>68.7</u>	<u>77.3</u>	4.3	38.6	50.2	<u>4.7</u>	30.0	<u>36.3</u>	4.8	44.2	<u>54.1</u>	4.7

Table 14: 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											
<i>Top-5 Documents</i>											
Original Docs	-	36.9	48.8	61.5	70.3	29.6	39.9	22.2	29.1	<u>37.5</u>	<u>47.0</u>
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	<u>61.0</u>	<u>70.0</u>	29.6	39.9	23.4	30.0	37.1	46.5
RECOMP-Extr	Ext.	<u>35.9</u>	45.8	<u>58.5</u>	<u>66.1</u>	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	<u>47.4</u>	<u>61.0</u>	69.8	<u>29.7</u>	<u>40.3</u>	<u>25.9</u>	<u>32.0</u>	38.1	47.4
<i>Top-20 Documents</i>											
Original Docs	-	39.2	51.4	64.2	73.2	30.6	40.8	24.8	<u>32.4</u>	<u>39.7</u>	<u>49.4</u>
RECOMP-Abst	Abs.	30.2	40.7	59.3	67.6	27.4	37.7	26.8	<u>32.4</u>	35.9	44.6
CompAct	Abs.	35.6	47.0	60.1	69.2	30.7	40.6	21.0	<u>28.3</u>	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	<u>63.8</u>	72.1	<u>31.0</u>	<u>41.1</u>	25.3	31.6	38.5	47.5
EXIT (Ours)	Ext.	<u>38.8</u>	<u>50.8</u>	63.4	<u>72.3</u>	32.2	43.1	<u>26.6</u>	32.9	40.3	49.8
Llama-3.1-70B-Instruct											
<i>Top-5 Documents</i>											
Original Docs	-	<u>39.5</u>	<u>51.8</u>	68.3	76.5	36.0	46.7	31.4	37.5	<u>43.8</u>	<u>53.1</u>
RECOMP-Abst	Abs.	<u>38.1</u>	<u>50.3</u>	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	<u>36.8</u>	<u>47.5</u>	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	<u>69.3</u>	<u>77.2</u>	35.6	46.8	<u>34.7</u>	<u>40.2</u>	43.7	52.8
EXIT (Ours)	Ext.	<u>39.5</u>	51.9	69.6	77.8	38.1	49.4	35.4	41.0	45.7	55.1
<i>Top-20 Documents</i>											
Original Docs	-	<u>42.1</u>	<u>55.0</u>	71.1	79.1	<u>39.4</u>	51.1	<u>35.4</u>	42.4	<u>47.0</u>	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	<u>35.4</u>	41.1	45.2	<u>54.6</u>
EXIT (Ours)	Ext.	42.5	55.3	<u>71.0</u>	<u>79.0</u>	39.8	51.1	36.5	<u>42.2</u>	47.5	56.9

Table 15: 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 (%) ↓
<i>Top-5 Documents</i>						
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.	38.0 (5.2)	36.9 (5.0)	63.3 (8.4)	55.1 (7.0)	46.0 (6.2)
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)	<u>76.6 (10.2)</u>	<u>62.1 (7.9)</u>	98.5 (13.4)
RECOMP-Extr	Ext.	<u>43.9 (6.1)</u>	<u>42.7 (5.8)</u>	90.2 (12.0)	97.9 (12.5)	<u>58.9 (8.0)</u>
LongLLMLingua	Ext.	<u>224.3 (31)</u>	<u>221.9 (30.4)</u>	229.2 (30.6)	234.5 (29.9)	<u>225.1 (30.7)</u>
Ours (EXIT)	Ext.	283.8 (39.2)	211.3 (28.9)	190.3 (25.4)	154.9 (19.7)	228.4 (31.2)
<i>Top-20 Documents</i>						
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.	26.1 (0.9)	38.6 (1.3)	64.5 (2.2)	51.1 (1.6)	43.1 (1.5)
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)	<u>92.7 (3.0)</u>	175.2 (6.0)
RECOMP-Extr	Ext.	<u>43.4 (1.5)</u>	<u>40.6 (1.4)</u>	89.3 (3.0)	<u>95.7 (3.0)</u>	<u>57.8 (2.0)</u>
LongLLMLingua	Ext.	<u>550.9 (19.0)</u>	<u>553.9 (18.9)</u>	563.3 (18.8)	595.5 (19.0)	<u>556 (18.9)</u>
Ours (EXIT)	Ext.	1001.2 (34.6)	635.0 (21.7)	465.0 (15.5)	367.1 (11.7)	700.4 (23.9)

Table 16: 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																
<i>Top-5 Documents</i>																
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	0.43	1.55	1.06	0.29	1.35	1.92	0.32	2.24	1.73	0.38	2.11	1.46	0.36	1.81
CompAct	Abs.	8.04	0.43	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.	0.04	<u>0.50</u>	0.54	0.04	<u>0.31</u>	0.35	0.04	<u>0.37</u>	0.41	0.04	0.59	0.63	0.04	0.44	0.48
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>	0.43	<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>
<i>Top-20 Documents</i>																
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	0.35	2.44	2.16	0.44	2.60	1.72	0.43	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.	0.11	0.51	0.62	0.11	0.28	0.39	0.12	0.42	0.54	0.11	0.56	0.68	0.11	<u>0.44</u>	0.55
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	0.42	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>	0.43	<u>1.85</u>
Llama-3.1-70B-Instruct																
<i>Top-5 Documents</i>																
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	2.14	3.34	2.06	2.37	4.43	1.71	2.54	4.24	1.56	2.56	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.	0.04	2.44	2.48	0.04	<u>2.21</u>	2.25	0.04	2.87	2.91	0.05	3.29	3.33	0.04	<u>2.70</u>	2.74
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>
<i>Top-20 Documents</i>																
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	2.75	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	1.50	48.90	31.40	1.50	33.00	9.70	3.10	12.80	30.25	2.90	33.15
RECOMP-Extr	Ext.	0.11	2.27	2.38	0.12	<u>2.10</u>	2.21	0.12	2.77	2.89	0.13	3.18	3.31	0.12	2.58	2.70
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