# FAMILIARITY-AWARE EVIDENCE COMPRESSION FOR RETRIEVAL-AUGMENTED GENERATION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Retrieval-augmented generation (RAG) improves large language models (LMs) by incorporating non-parametric knowledge through evidence retrieved from external sources. However, it often struggles to cope with inconsistent and irrelevant information that can distract the LM from its tasks, especially when multiple evidence pieces are required. While compressing the retrieved evidence with a compression model aims to address this issue, the compressed evidence may still be unfamiliar to the target model used for downstream tasks, potentially failing to utilize the evidence effectively. We propose FAVICOMP (FAmiliarity-aware EVIdence COMPression), a novel training-free evidence compression technique that makes retrieved evidence more familiar to the target model, while seamlessly integrating parametric knowledge from the model. Specifically, FAVICOMP proactively composes the compressed evidence in a way to lower the perplexity of the target model by combining decoding probabilities from both the compression model and the target model to generate context that is more familiar to the target model. This approach balances the integration of parametric and non-parametric knowledge, which is especially helpful in complex tasks where the retrieved evidence set may not contain all the necessary information. Experimental results show that FAVICOMP consistently outperforms most recent evidence compression baselines across multiple open-domain QA datasets, improving accuracy by up to 23.91% while achieving high compression rates. Additionally, we demonstrate the effective integration of both parametric and non-parametric knowledge during evidence compression.

031 032 033

034

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

#### 1 INTRODUCTION

Retrieval-augmented generation (RAG) has become a common paradigm for large language models
(LMs) to leverage external knowledge beyond their inherent knowledge boundaries to perform better
in knowledge-intensive tasks such as open-domain question answering (QA) (Lewis et al., 2020;
Izacard & Grave, 2021; Guu et al., 2020) and fact-checking (Pan et al., 2023; Li et al., 2024c). In
particular, incorporating multiple evidence pieces is crucial in solving complicated tasks such as
multi-hop and complex reasoning (Trivedi et al., 2023; Jiang et al., 2023b; Li et al., 2024b; Lu et al., 2023), which require various sources of information to solve the questions.

Nevertheless, RAG often struggles to cope with inconsistent and irrelevant information from the 043 multiple evidence set, which can interfere with downstream tasks (Shi et al., 2023). This highlights 044 the need for *compression-based* RAG (Jiang et al., 2023a; Xu et al., 2024; Yoon et al., 2024) to identify and retain only the essential information for the LMs to utilize effectively. Traditionally, 046 compression-based RAG has focused on reranking documents or sentences by relevance and then 047 incorporating a top-ranked subset (Nogueira et al., 2020; Zhuang et al., 2023; Wang et al., 2023c) 048 or compressing the documents into an abstractive summary that retains only essential context (Jiang et al., 2023a; Xu et al., 2024; Yoon et al., 2024). However, the compressed evidence might be unfamiliar to the LM employed for the downstream task (referred to as the target model), particularly 051 due to discrepancies in the pretrained internal knowledge and prompt preferences between the compression model and the target model (Gonen et al., 2023; Lee et al., 2024; Li et al., 2024a; Mallen 052 et al., 2023). When LMs encounter unfamiliar contextual information, they often fail in balancing parametric and non-parametric knowledge, either by overly relying on their parametric knowledge



Figure 1: An overview of FAVICOMP. Instead of relying solely on compressed evidence from the compression model (upper), FAVICOMP familiarizes the compressed evidence to the target model while integrating parametric knowledge through ensemble decoding, resulting in improved downstream performance (lower).

(Longpre et al., 2021; Wang et al., 2023a; Zhou et al., 2023) or by utilizing retrieved evidence without considering its relevance to the input (Wu et al., 2024).

073 To address these challenges, we propose FAVICOMP (FAmiliarity-aware EVIdence COMPression), a training-free evidence compression method that makes retrieved multi-evidence more familiar to 074 the target model, while seamlessly integrating parametric knowledge from the model. Inspired by 075 the prior findings that an LM's familiarity with a prompt is generally reflected by low perplexity 076 (Liu et al., 2024; Gonen et al., 2023; Wang et al., 2023b), FAVICOMP proactively composes the 077 compressed evidence in a way to lower the perplexity of the target model. Specifically, FAVICOMP leverages the decoding probabilities of two LMs, a *compression model* and the *target model*. The 079 compression model is instructed to summarize the raw evidential documents into a relevant context to the input, while the target model is instructed to generate relevant context without referencing the 081 documents. Instead of directly selecting the highest probability token from the compression model at 082 each decoding step, we ensemble the token logits from both the compression and target models and 083 then select the token with the highest probability from this combined set. This ensemble decoding 084 therefore constrains the token search space of the compression model to those with lower perplexity for the target model, making the context more familiar to the target model (Liu et al., 2024). 085

Furthermore, FAVICOMP potentially synergizes the retrieved knowledge with the target model's parametric knowledge introduced during ensemble decoding. FAVICOMP can effectively discern when to leverage internal or external knowledge, which is particularly beneficial in the presence of noisy contextual evidence in complex tasks such as multi-document or multi-hop QA (Wang et al., 2024).

FAVICOMP brings along key advantages of RAG for complex tasks from two perspectives. On the one hand, it is capable of compressing multiple augmented documents to a more favorable form to the target model. This mechanism not only helps the model better comprehend the essential evidence in the retrieval augmentation but also better balances knowledge utility in both the evidential context and the model's parametric memory. On the other hand, it is a training-free and model-agnostic approach that can be easily plugged into any RAG processes

Our experiments show that FAVICOMP outperforms most recent evidence compression baselines in five open-domain QA datasets, improving accuracy by up to 23.91% while maintaining high compression rates. Additionally, we conduct ablation studies by varying the degree of decoding ensemble and analyzing its impact on performance and context perplexity. Moreover, we investigate how FAVICOMP effectively integrates parametric and non-parametric knowledge during evidence compression.

103 104

066

067

068

069

#### 2 Method

105 106

107 We present FAVICOMP, a decoding-time evidence compression method that familiarizes retrieved evidence with the target model while synergizing them with the model's parametric knowledge. We

first illustrate the motivation for FAVICOMP in §2.1 and provide the preliminaries of compression based RAG in §2.2, followed by a detailed definition of our proposed framework in §2.3.

111 2.1 MOTIVATION AND METHOD OVERVIEW

113 Standard RAG faces the challenge of LMs struggling to address inconsistent and irrelevant information from multiple evidence pieces, which can interfere with downstream tasks (Shi et al., 114 2023). Previous research has primarily concentrated on question-focused compression (Jiang et al., 115 2023a; Xu et al., 2024; Yoon et al., 2024); however, this approach may lead to suboptimal perfor-116 mance in downstream tasks due to the compressed evidence's potential unfamiliarity with the target 117 model employed. This unfamiliarity arises from discrepancies in pretrained internal knowledge and 118 prompt preferences between the compression model and the target model (Gonen et al., 2023; Lee 119 et al., 2024; Mallen et al., 2023). Furthermore, the unfamiliarity often leads to failure in balancing 120 parametric and non-parametric knowledge, either by overly relying on their parametric knowledge 121 (Longpre et al., 2021; Wang et al., 2023a; Zhou et al., 2023) or by using retrieved evidence without 122 considering its relevance to the input (Wu et al., 2024). To address this issue, FAVICOMP introduces 123 a novel approach that compresses evidence that better aligns with the target model's preferences 124 while seamlessly integrating parametric knowledge into the compressed evidence using a novel en-125 semble decoding technique, thereby improving its performance on downstream tasks.

Fig. 1 illustrates the overview of FAVICOMP. In this example, FAVICOMP makes the compressed evidence more favorable to the target model and leverages its parametric knowledge to supplement the missing evidence (*"Lionel Messi made his league debut in Barcelona"*), effectively combining evidential and parametric knowledge.

130 131

132

140 141 142

152

153

154

156 157

158 159

160 161

#### 2.2 COMPRESSION-BASED RETRIEVAL AUGMENTED GENERATION

Given a set of k retrieved evidence snippets  $D = \{d_1, d_2, \dots, d_k\}$  and a textual input sequence x, standard RAG aims to generate an output sequence y, conditioned on both D and x. However, standard RAG directly utilizes D which often contains irrelevant information to x, potentially confusing the target model in downstream tasks (Shi et al., 2023). Thus, the compression-based RAG uses an additional compression model to condense D into a concise and input-relevant context c, which is then used in place of D during the downstream generation process. Thus, the compression-based RAG is formalized as:

$y^* = \arg\max_{y} P_t(y \mid x, \hat{c}),$
$\hat{c} = P_c(c \mid x, [d_1, d_2, \dots, d_k]),$

where  $y^*$  is the final output sequence,  $[\cdot, \cdot]$  denotes concatenation, and  $P_t$  and  $P_c$  represent the probability distributions of the target and compression models, respectively. In this work, we consider any natural language prompting tasks, such as open-domain QA tasks, where x represents the input prompt (also known as the query in QA tasks) and  $y^*$  denotes the output sequence.

The compression model's objective is to produce a concise yet informative summary c of the evidential documents D that captures the essential information relevant to the input query x. We use an unsupervised approach, where the model is instructed to generate a query-relevant summary of D in a zero-shot manner using an evidence compression instruction prompt, denoted as  $I_{comp}$ , such as the one below:

**Evidence Compression Instruction** 

Given a question and multiple document snippets, generate one summarized context that is helpful to answer the question.

Specifically, the evidence compression is done in an auto-regressive way formalized as,

$$P_{c}(c \mid I_{comp}, x, D) = \prod_{i=1}^{|c|} P_{c}(c_{i} \mid I_{comp}, x, D, c_{< i}),$$

where |c| is the length of the summary c.

## 162 2.3 ENSEMBLE DECODING FOR FAVICOMP

Simple compression techniques might lead to subpar performance in downstream tasks because the compressed evidence may not be familiar to the target model. To better align the context to the target model, FAVICOMP proactively composes it to lower the target model's perplexity by introducing a constraint in decoding space from the target model during the evidence compression. FAVICOMP achieves this goal through ensemble decoding, which involves a multiplicative ensemble of two LMs—compression model and target model—at each decoding step.

Specifically, the target model is directed to generate a context c that would be helpful in answering the question x without referencing the evidence set. This is also done in zero-shot using a context generation instruction prompt  $I_{gen}$  such as:

173 174 175

176 177 178

179

181 182

189 190 Context Generation Instruction

Given a question, generate a context that is helpful to answer the question.

The context generation is also performed in an auto-regressive fashion, represented as:

$$P_t(c \mid I_{gen}, x) = \prod_{i=1}^{|c|} P_t(c_i | I_{gen}, x, c_{< i}),$$

where |c| denotes the length of the generated context c.

Once the compression model and the target model generate their respective probability distributions for the next token, the subsequent token is chosen by maximizing the weighted sum of the log probabilities from both models. The selected token is the continuation of the previously generated text aligned with their objectives. This process is formalized as follows:

$$c_i = \arg\max_{c_i \in V} ((1 - \alpha) \cdot \log P_c(c_i | I_{comp}, x, D, c_{< i}) + \alpha \cdot \log P_t(c_i | I_{gen}, x, c_{< i})),$$

where  $c_i$  is the subsequent token, and  $\alpha$  is the ensemble coefficient that weighs between the two probability distributions. We demonstrate how the coefficient  $\alpha$  impacts both the perplexity and the downstream performance in §4.2.

Ensemble decoding proactively shifts the token search space in evidence compression by upweighting those tokens with lower perplexity from the target model's perspective (Liu et al., 2024), resulting in a compressed evidence that is more familiar to the target model. Note that since both objectives ultimately share the goal of generating context relevant to the question, combining the logits ensures alignment with this ultimate goal.

In addition, ensemble decoding enables FAVICOMP to seamlessly integrate both retrieval knowl-200 edge from the external evidence set and the target model's parametric knowledge. Specifically, 201 FAVICOMP selects the arg max token from the target model only when the token's probability is 202 higher than that of the compression model, demonstrating that FAVICOMP draws on parametric 203 knowledge only when necessary-potentially when the compression model is uncertain about the 204 next token. This is particularly beneficial for complex tasks like multi-document QA, where the 205 evidence set may not include all the necessary information (Mallen et al., 2023). In such cases, the 206 missing information in compressed evidence can be supplemented by tokens generated from context 207 generation by the target model, which is entirely based on parametric knowledge. We demonstrate 208 in §4.3 and §5 that FAVICOMP can incorporate knowledge from both sources effectively, leading to 209 a performance boost compared to compression methods that solely focus on distilling knowledge 210 from the evidence set.

211

#### 212 213 3 EXPERIMENTAL SETTINGS

214

# We assess the effectiveness of FAVICOMP on knowledge-intensive QA tasks. In this section, we delve into the details of the experimental settings.

### 216 3.1 DATASETS

We evaluate FAVICOMP on five open-domain QA datasets, including two single-document QA datasets, Natural Questions (NQ) (Kwiatkowski et al., 2019) and TriviaQA (TQA; Joshi et al. 2017), and three multi-document QA datasets, HotpotQA (Yang et al., 2018), 2WikiMultiHopQA (2WikiMQA; Ho et al. 2020), and MuSiQue (Trivedi et al., 2022). Following prior studies (Asai et al., 2023; Xu et al., 2024), we evaluate the performance on the development set of each dataset and use three evaluation metrics, i.e. Accuracy (Acc), token-level F1 and compression rate (Comp) which is calculated as  $\frac{\# of tokens in retrieved documents}{\# of tokens in compressed documents}$ .

225 226

239 240

241

#### 3.2 IMPLEMENTATION DETAILS

For all the comparison methods, we utilize three LMs as the target model to tackle down-stream QA tasks with RAG, i.e. Llama3-8B-Instruct<sup>1</sup>, Mistral-7B-Instruct<sup>2</sup> and Mixtral-8x7B-Instruct<sup>3</sup>. For each question, we retrieve five documents from 2018 Wikipedia corpus (Karpukhin et al., 2020) using Contriever-MSMARCO<sup>4</sup> (Izacard et al., 2021), so as to be consistent with previous studies (Xu et al., 2024; Yoon et al., 2024).

232 For FAVICOMP, we employ three compression and target model pairs: (1)233 Llama3.2-3B-Instruct and Llama3-8B-Instruct as the compression models 234 and Llama3-8B-Instruct as the target model, (2) Mistral-7B-Instruct as the 235 compression model and Mixtral-8x7B-Instruct as the target model, and (3) same 236 Mistral-7B-Instruct as the compression model and the target model (Appx. §B.1). Also, 237 we set  $\alpha$  to 0.5 by default, for which more analyses are given in §4.2. The prompts used in the 238 experiment are presented in Appx. §C.

#### 3.3 BASELINES

We consider the following categories of baselines. (1) No Context: RAG without any context. (2) 242 Gold Compression: RAG using directly relevant evidence from the retrieved documents if they 243 exist. (3) **Raw Document**: RAG with raw documents that have not undergone any compression. (4) 244 Generated Context (Yu et al., 2023): RAG with context generated by the same LM as the target 245 model. This is equivalent to FAVICOMP with  $\alpha = 1$ , as we rely solely on the target model to generate 246 context when  $\alpha = 1$ . (5) **Reranking-based Methods**: We rerank sentences in the evidence set and 247 choose top-ranked sentences as the context. We utilize two rerankers-Sentence-BERT (Reimers & 248 Gurevych, 2020) and RECOMP-extractive (Xu et al., 2024). (6) Compression-based Methods: We 249 employ four compressors—LongLLMLingua (Jiang et al., 2023a), RECOMP-abstractive (Xu et al., 250 2024), CompAct (Yoon et al., 2024), and Zero-shot Summarization. Zero-shot Summarization is 251 instructed to summarize the evidence set into a concise summary based on the question, using the 252 same LM as the target model. This is equivalent to FAVICOMP with  $\alpha = 0$ , as we depend entirely 253 on the compression model without any intervention from the target model. A detailed explanation of the implementation of the baselines is provided in Appx. §A. 254

255 256

257

262

263 264

265

266

#### 4 EXPERIMENTAL RESULTS

In this section, we compare the overall performance of FAVICOMP with other baselines across the five datasets (§4.1), explore the impact of ensemble coefficient  $\alpha$  on performance and perplexity (§4.2), investigate how effectively FAVICOMP incorporate parametric and non-parametric knowledge (§4.3), and compare the compression rates with other baselines (§4.4).

4.1 MAIN RESULTS

The overall performance of FAVICOMP and the baselines across the five datasets are presented in Tab. 1 and Tab. 3. To start with, the compression-based methods consistently outperform the

<sup>267 &</sup>lt;sup>1</sup>https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct

<sup>268 &</sup>lt;sup>2</sup>https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1

<sup>&</sup>lt;sup>4</sup>https://huggingface.co/facebook/contriever-msmarco

Methods	Size		NQ			TQA		1	Hotpot	QA	2	WikiM	QA	1	MuSiQ	ue
Wiemous	5120	Acc	F1	Comp	Acc	F1	Comp	Acc	F1	Comp	Acc	F1	Comp	Acc	F1	Comp
Llama3-8B-Instruct																
Gold Compression	-	-	-	-	-	-	-	42.3	51.3	-	35.7	40.0	-	10.2	17.7	-
No Context Raw Document Generated Context	- -	26.9 42.6 32.3	31.9 <b>47.1</b> 36.6	- -	57.2 67.6 59.7	61.2 70.8 62.4	- - -	19.1 30.3 22.7	25.5 38.7 29.7	- -	20.5 22.0 24.8	25.0 26.8 28.7	- - -	5.4 8.2 7.6	13.0 15.0 14.8	- -
Sentence-BERT RECOMP-extractive	110M 110M <sup>†</sup>	30.3 33.7	35.4 38.1	21.13 19.45	59.2 59.4	62.9 62.8	20.61 18.86	22.4 22.5	29.6 29.8	10.30 9.47	18.1 18.0	22.9 22.4	9.96 9.17	7.7 8.1	14.8 15.5	10.18 9.24
LongLLMLingua RECOMP-abstractive CompAct	$\begin{array}{c} 7B^\dagger \\ 775M^\dagger \\ 7B^\dagger \end{array}$	35.4 39.3 42.3	40.9 43.3 46.1	1.87 17.96 8.85	64.8 62.9 67.0	67.6 66.1 69.7	1.84 17.79 8.92	25.9 27.0 29.8	34.7 34.8 37.5	1.83 19.72 9.45	19.2 20.5 21.4	24.2 25.0 26.6	1.83 32.06 10.71	7.7 7.3 9.2	14.4 14.8 16.9	1.83 32.05 8.96
Zero-shot Summarization	3B 8B	39.4 41.3	43.2 45.1	14.12 13.87	64.2 66.3	67.1 69.5	17.12 16.58	30.1 30.2	38.5 38.6	18.75 17.38	25.7 22.3	31.1 28.1	21.39 18.98	7.7 8.3	15.3 16.3	16.19 15.50
FAVICOMP	3B 8B	<b>42.8</b> 42.3	46.8 46.6	16.43 15.79	68.0 68.4	70.9 <b>71.5</b>	22.40 20.99	<b>33.0</b> 32.3	<b>41.6</b> 41.0	22.55 21.49	<b>29.6</b> 27.6	<b>35.2</b> 33.6	23.10 22.41	10.8 11.4	19.9 <b>20.1</b>	18.95 19.06
					Λ	1ixtral-	8x7B-Inst	ruct								
Gold Compression	-	-	-	-	-	-	-	48.2	55.1	-	49.9	51.9	-	12.9	18.6	-
No Context Raw Document Generated Context	- -	36.7 <b>46.3</b> 33.6	38.4 42.1 33.9	- -	68.9 72.1 61.4	72.0 71.1 62.9	- - -	25.1 34.0 26.5	31.6 39.0 32.9	- - -	32.5 32.9 30.2	35.9 36.3 34.3	- -	6.4 10.1 7.2	11.8 15.6 13.4	- - -
Sentence-BERT RECOMP-extractive	110M 110M <sup>†</sup>	36.8 38.0	36.8 37.9	21.13 19.42	67.0 66.7	68.7 68.0	20.61 18.81	28.3 28.7	34.5 34.3	10.13 9.30	32.5 31.8	36.2 34.9	9.76 9.01	9.9 9.4	15.2 15.6	10.07 9.11
LongLLMLingua RECOMP-abstractive CompAct	$\begin{array}{c} 7B^{\dagger} \\ 770M^{\dagger} \\ 7B^{\dagger} \end{array}$	40.1 42.1 44.1	39.4 41.3 43.4	1.96 17.55 8.83	70.5 68.4 70.3	71.0 69.4 71.4	1.96 17.47 8.92	32.0 32.3 35.2	38.3 38.5 41.6	1.95 19.39 9.45	31.9 32.2 35.9	36.1 36.2 39.5	1.93 31.20 10.67	9.7 7.9 11.2	15.9 13.6 16.9	1.96 31.18 8.94
Zero-shot Summarization FAVICOMP	7B 7B	42.1 43.6	40.6 <b>44.5</b>	8.65 7.30	65.9 7 <b>2.6</b>	67.0 <b>73.9</b>	10.43 8.21	31.4 36.3	38.1 <b>44.4</b>	11.71 8.89	28.5 <b>40.5</b>	32.8 <b>45.2</b>	14.35 10.26	8.4 13.4	13.8 <b>19.9</b>	10.26 8.42

Experimental results on five open-domain QA datasets. 291 Table 1: Size column repre-<sup>†</sup> indicates a fullysents the size of the compression model used for each method. 292 supervised compression model, where the reranker or the compressor is trained. For 293 the experiment with Llama3-8B-Instruct, Zero-shot Summarization and FAVICOMP use Llama3.2-3B-Instruct and Llama3-8B-Instruct as the compression model, shown as 295 3B and 8B in the Size column. The best Accuracy and token-level F1 scores for each dataset are in 296 bold. 297

298

299

300

reranking-based methods, due to the fact the reranking-based methods are prone to losing more 301 question-relevant information by discarding lower-ranked sentences. Next, FAVICOMP outperforms 302 all other baselines across all the datasets, except for the Gold Compression which is regarded as the 303 upper bound of the performance. It is noteworthy that FAVICOMP, as a training-free, decoding-time 304 strategy, outperforms supervised baselines even with the 3B parameters compression model. For the 305 MuSiQue dataset, FAVICOMP even outperforms Gold Compression baseline which can be viewed 306 as a perfect compressor. This demonstrates that explicitly incorporating parametric knowledge from 307 the target model can significantly enhance performance in multi-document QA, even when the con-308 text is imperfect.

309 Moreover, it is surprising that most of the supervised compression-based methods are excelled by 310 the Raw Document baseline. This indicates that existing methods are likely to fall short of re-311 taining essential supportive information while compressing the evidence documents. Additionally, 312 LongLLMLingua and RECOMP-abstractive perform worse than Zero-shot Summarization with 313 similar or smaller size compression model. This may be possibly due to the use of smaller base 314 model for compression (T5-large for RECOMP-abstractive), but it also suggests that knowledge 315 distillation from larger teacher LM to the smaller compression model may not generalize well, as the context preferences and prior knowledge of the target model and the teacher model are likely 316 to differ. We conduct a head-to-head experiment on RECOMP-abstractive by using the same base 317 compression model as FAVICOMP for a more fair comparison in Appx. §B.2. 318

Furthermore, despite using the same base model for the compression model
 (Mistral-7B-Instruct), the training-free FAVICOMP outperforms CompAct, which
 trains the compression model using knowledge distillation to generate and evaluate summaries of
 retrieved documents. This also indicates that knowledge distilled from a teacher model may not
 always be effectively transferable to the target model due to discrepancies in context preference and
 prior knowledge. In contrast, the superior performance of FAVICOMP is attributed to its ability to



Figure 2: Impact of coefficient  $\alpha$  on performance and perplexity.

familiarize evidence with the target model and its effective incorporation of parametric knowledge from ensemble decoding.

Finally, given that Zero-shot Summarization corresponds to FAVICOMP with  $\alpha = 0$  and Generated Context corresponds to FAVICOMP with  $\alpha = 1$ , the fact that FAVICOMP outperforms both baselines highlights its ability to effectively incorporate tokens from both sources—evidence summary and generated context. This results in superior performance compared to relying on just one source alone.

344 345 346

334

335 336 337

338

#### 4.2 IMPACT OF ENSEMBLE COEFFICIENT ON PERFORMANCE AND PERPLEXITY

347 Fig. 2 illustrates how performance and perplexity change as the ensemble coefficient  $\alpha$  is varied 348 across the values  $\{0.0, 0.1, 0.3, 0.5, 0.7, 0.9, 1.0\}$  on NQ, HotpotQA and MuSiQue datasets. We 349 calculate the perplexity of the compressed evidence conditioned on the preceding inputs, i.e. in-350 struction, demonstrations, and the question. For all the datasets, performance is the highest when 351  $\alpha = 0.5$ , indicating that proactively lowering perplexity by equally weighting both input sources yields the best results. When  $\alpha$  is below 0.5, performance improves as the perplexity of compressed 352 evidence decreases, which aligns with the previous works (Liu et al., 2024; Gonen et al., 2023). 353 However, when  $\alpha$  exceeds 0.5, performance declines as perplexity decreases due to the lack of ev-354 idential knowledge during evidence compression. Additionally, when  $\alpha$  reaches 0.9 or 1.0, there is 355 a slight rise in the perplexity due to LM's increased uncertainty with limited evidential knowledge. 356 Results for other datasets are included in Fig. 4.

357 358 359

360

#### 4.3 INTEGRATION OF PARAMETRIC AND NON-PARAMETRIC KNOWLEDGE

The effective integration of parametric and non-parametric knowledge is crucial for complex tasks 361 such as multi-document QA, where the evidence set may not contain all the necessary information. 362 To this end, we evaluate how effectively FAVICOMP incorporates parametric knowledge from the 363 target model and non-parametric knowledge from the compression model on the multi-document 364 QA datasets. We begin by dividing the test samples of each dataset into evidence-relevant and evidence-irrelevant subsets, using the *Hits* metric. The *Hits* metric is set to 1 (evidence-relevant) 366 if the retrieved evidence set contains the correct answer, and 0 (evidence-irrelevant) if it does not. We 367 then assess the downstream performance of each subset. The underlying intuition is that if a method 368 performs better on the evidence-relevant subset, it suggests that the method is more effectively uti-369 lizing the provided evidential knowledge. Conversely, if a method excels on the evidence-irrelevant subset, it indicates that the method is more effectively leveraging parametric knowledge without 370 relying on potentially irrelevant evidence. 371

The left figure of Fig. 3 compares the accuracy in Hits = 0 and Hits = 1 subsets across the datasets. We compare FAVICOMP with the top-performing unsupervised compression method, Zeroshot Summarization, and the most competitive supervised compression method, CompAct. Compared to the other two baselines, FAVICOMP performs better in the Hits = 0 subset while performing comparably in the Hits = 1 subset. This proves that FAVICOMP effectively relies on parametric knowledge rather than evidential knowledge when faced with irrelevant evidence, while maintaining similar effectiveness in utilizing evidential knowledge when relevant evidence is present.



Figure 3: Accuracy of baselines (left) and FAVICOMP with various  $\alpha$  values (right) on Hits = 0and Hits = 1 subset of multi-document QA datasets.

Interestingly, even though CompAct generally performs better on the Hits = 1 subset compared to Zero-shot Summarization, it underperforms relative to Zero-shot Summarization on the Hits = 0 subset. This suggests that the training may have been biased towards utilizing solely evidential knowledge, rather than effectively leveraging both sources in synergy.

We also evaluate the performance of FAVICOMP with various  $\alpha$  values under this setting. The right figure of Fig. 3 shows that  $\alpha = 0.5$  or  $\alpha = 0.7$  performs the best on the Hits = 0 subset, while performance declines as  $\alpha$  deviates further from the value. This pattern in the Hits = 0 subset mirrors the overall performance trend, suggesting that appropriately utilizing parametric knowledge when the evidence is irrelevant is crucial to the overall performance. In the Hits = 1 subset, performance remains consistent for  $\alpha$  values up to 0.5 but decreases significantly when  $\alpha$  exceeds 0.5 due to the diminished utilization of the relevant evidential context.

404 405

388

389

390 391 392

393

394

395

396

#### 4.4 COMPRESSION RATE COMPARISONS

406 Since one of the functionalities of compression-based RAG is to reduce the number of tokens from 407 the evidence while keeping its essential information, we report the compression rate in Tab. 1. Over-408 all, reranking-based methods, RECOMP-abstractive and FAVICOMP consistently score the highest 409 compression rates. Reranking-based methods achieves high compression since they only select one 410 or two sentences that may contain the answer to the question, but the information loss is more signif-411 icant compared to other methods. RECOMP-abstractive exhibits high compression rates because the 412 compression model is trained to output an empty string when no relevant evidence is found, which is often the case in multi-document QA datasets. FAVICOMP compresses the evidence to make it 413 familiar to the target model by lowering its perplexity at each decoding step, typically resulting in 414 a shorter context. Notably, when compared to Zero-shot Summarization, which is equivalent to 415 FAVICOMP with  $\alpha = 0$ , FAVICOMP consistently achieves higher compression rates. This demon-416 strates that the ensemble decoding strategy, combining token logits from both evidence compression 417 and context generation, leads to greater compression efficiency. 418

419

#### 5 CASE STUDY

420 421

Tab. 2 presents two examples from HotpotQA to illustrate how FAVICOMP effectively familiarizes evidence while seamlessly integrating both parametric and non-parametric knowledge during evidence compression. We compare its output with Raw Document, which does not apply any compression, and Zero-shot Summarization, which is equivalent to FAVICOMP with  $\alpha = 0$ .

In both examples, Raw Document fails to produce the correct answer, even though the evidence
contains the necessary information, highlighting the need for effective evidence compression. In the
first example, while the difference between the compressed evidence from Zero-shot Summarization and FAVICOMP appears subtle, FAVICOMP delivers the correct answer with a lower perplexity
in compression, underscoring the significance of evidence familiarization. The second example
highlights the importance of parametric knowledge when the retrieved evidence set lacks complete
information. Since the evidence set does not mention "Skeptic," Zero-shot Summarization intro-

Methods	(Compressed) Evidence	Prediction	Perplexity
Raw Document	( <i>skip</i> ) The Tales of Hoffmann is a 1951 British Tech- nicolor film adaptation of Jacques Offenbach's opera "The Tales of Hoffmann", written, produced and directed by the team of Michael Powell and Emeric Pressburger working un- der the umbrella of their production company, The Archers.	Emeric Pressburger 🗡	12.429
Zero-shot Summarization	The 1951 film "The Tales of Hoffmann" is an adaptation of Jacques Offenbach's opera, written, produced, and directed by Michael Powell and Emeric Pressburger.	Emeric Pressburger X	2.298
FaviComp	The 1951 film "The Tales of Hoffmann" is an adapta- tion of Jacques Offenbach's opera, written by Emeric Press- burger, a Hungarian-British screenwriter, and directed by Michael Powell and Emeric Pressburger.	The Tales of Hoffmann ✓	1.959
Quest	tion: Which magazine was first published earlier, The Chronicle	of Philanthropy or Skepti	c?
Methods	(Compressed) Evidence	Prediction	Perplexit
Raw Document	The Chronicle of Philanthropy is a magazine that covers the nonprofit world( <i>skip</i> ) It was founded in 1988 by editor Phil Semas and then managing editor Stacy Palmer. ( <i>skip</i> ) Philanthropy (magazine) Philanthropy is a quar- terly magazine published by the Philanthropy Roundtable. First published as a newsletter in 1987, "Philanthropy" be- came a glossy magazine in 1996.	Philanthropy 🗡	4.856
Zero-shot Summarization	The Chronicle of Philanthropy was founded in 1988, while Philanthropy magazine was first published as a newsletter in 1987 and became a glossy magazine in 1996.	Philanthropy magazine X	3.196
	The Chronicle of Dhilenthrony was first published in 1099	The Chronicle of	L

Table 2: Case study of evidence compression: FAVICOMP vs. Raw Document and Zero-shot Summarization. For FAVICOMP, the colors red and blue highlight tokens that are the arg max of the compression model and the target model, respectively. Purple indicates a token that is the arg max of neither model. Tokens with no coloring represent those that are the  $\arg \max$  of both models.

duces irrelevant information ("Philanthropy magazine"), ultimately leading to an incorrect answer. In contrast, FAVICOMP integrates parametric knowledge about "Skeptic" and incorporates it into the evidence compression. Notably, FAVICOMP selects the arg max token from the target model only when the token's probability is higher than that of the compression model, demonstrating that FAVICOMP draws on parametric knowledge only when necessary-potentially when the compres-sion model is uncertain about the next token. 

**RELATED WORKS** 

Evidence Compression for RAG. Standard RAG retrieves textual evidence related to the prompt from the external corpora or knowledge bases and incorporates it as a part of the input to the LM (Lewis et al., 2020; Izacard & Grave, 2021; Guu et al., 2020). However, retrieved evidence pieces may contain inconsistent or irrelevant information to the question, potentially confusing the target model in downstream tasks (Shi et al., 2023). To tackle this problem, traditional approaches aim to rerank the textual evidence based on its relevance to the question and then select a top-ranked subset to include as part of the input to the LM (Nogueira et al., 2020; Zhuang et al., 2023). However, this approach loses more question-relevant information by discarding lower-ranked sentences. 

Recent efforts on evidence compression seek to compress retrieved evidence pieces to filter out un-necessary information and retain only the essential context (Wang et al., 2023c; Li et al., 2024d; Ke et al., 2024; Jiang et al., 2023a; Xu et al., 2024; Cao et al., 2024; Yoon et al., 2024). Wang et al. (2023c) filter query-relevant context using relevance metrics and Li et al. (2024d) extract queryrelevant information and restructure them to form a consistent context. Ke et al. (2024) trains a seq2seq bridge model using supervised and reinforcement learning to optimize the connection be-tween the retriever and the LLM. Jiang et al. (2023a) and Cao et al. (2024) conduct token-level

486 or embedding-based compression to preserve only the query-relevant information using a trained 487 compressor. Xu et al. (2024) and Yoon et al. (2024) train a compression model to generate an ab-488 stractive summary of the documents by distilling knowledge from larger language models. While 489 these methods are successful to some extent, they often achieve suboptimal performance because 490 the compressed context may be unfamiliar to the LM used in the downstream task due to differences in pretrained internal knowledge and prompt preferences between the compression and the 491 target model. In contrast, FAVICOMP proactively compresses the evidence pieces in a way to lower 492 the target model's perplexity using an ensemble decoding technique without any training, thereby 493 improving the downstream performance. 494

495

Parametric and Non-parametric Knowledge in RAG. While there have been studies on the phe-496 nomena of LM's utilization of both parametric and non-parametric knowledge sources (Longpre 497 et al., 2021; Wadhwa et al., 2024; Wu et al., 2024; Zhang et al., 2024; Zhou et al., 2023; Wang et al., 498 2023a; Fang et al., 2024), there is a lack of research focused on effectively synergizing both sources. 499 A few of these efforts introduce counterfactual augmentation (Longpre et al., 2021; Fang et al., 2024; 500 Zhang et al., 2024) and causal intervention (Zhou et al., 2023; Wang et al., 2023a) to mitigate knowl-501 edge conflict, which, however, requires explicitly knowing the features of the input that causes such conflict. Zhang et al. (2023) seek to address this issue by incorporating LM-generated context into 502 the LM's input along with the retrieved documents, thereby integrating both sources of knowledge. 503 However, merely concatenating both contexts is a suboptimal solution, as LMs may still show bias 504 toward one source over the other when generating responses (Longpre et al., 2021; Wu et al., 2024). 505 To address this, FAVICOMP employs ensemble decoding during the evidence compression, ensuring 506 that both types of knowledge are seamlessly fused together to create a consistent context. 507

508 **Constrained Decoding.** Constrained decoding has been previously proposed in text generation 509 tasks for various purposes, including optimizing prompts (Liu et al., 2024), enhancing plausibility 510 (Li et al., 2023) or controllability (Meng et al., 2022; Huang et al., 2023), and reducing hallucination 511 (Shi et al., 2024). Contrastive Decoding (Li et al., 2023) enforces a plausibility constraint during generation by inducing the difference in token log-probabilities between expert and amateur LMs. 512 Context-aware Decoding (Shi et al., 2024) uses contrastive decoding to amplify the probability dif-513 ferences between outputs with and without evidence, encouraging the LM to prioritize the evidential 514 knowledge. Our work is closely connected with the method by Liu et al. (2024) which employs en-515 semble decoding to paraphrase prompts to enhance zero-shot LM prompting and generalization. 516 Their approach focuses on the robustness and generalizability of instruction prompts for tasks with-517 out retrieval augmentation. In contrast, our approach compresses externally retrieved evidence while 518 integrating parametric knowledge during compression, specifically targeting knowledge-intensive 519 tasks that require balancing both evidential and parametric knowledge.

520 521

#### 7 CONCLUSION

522 523 524

525

526

527

528

529

530

531

532

In this study, we introduce FAVICOMP, a training-free evidence compression method designed to enhance RAG by making retrieved evidence set more familiar to the target model, while seamlessly integrating parametric knowledge. By leveraging ensemble decoding, FAVICOMP compresses the retrieved evidence to make it more favorable to the target model. Moreover, FAVICOMP effectively balances the target model's parametric knowledge and the retrieved knowledge, improving performance on complex tasks where the retrieved evidence set may not contain all the necessary information. Our extensive experiments validate the effectiveness of FAVICOMP on open-domain QA tasks, showing significant improvements over recent evidence compression baselines in multiple datasets. Additionally, FAVICOMP's model-agnostic nature allows it to be effortlessly incorporated into various RAG workflows without additional training, making it a versatile tool for enhancing LMs in complex tasks.

533 534 535

#### 536 REFERENCES 537

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*, 2023.

540 Zhiwei Cao, Qian Cao, Yu Lu, Ningxin Peng, Luyang Huang, Shanbo Cheng, and Jinsong Su. 541 Retaining key information under high compression ratios: Query-guided compressor for LLMs. 542 In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics 543 (Volume 1: Long Papers), pp. 12685–12695, 2024. 544 Tianging Fang, Zhaowei Wang, Wenxuan Zhou, Hongming Zhang, Yanggiu Song, and Muhao Chen. 545 Getting sick after seeing a doctor? diagnosing and mitigating knowledge conflicts in event tem-546 poral reasoning. In Findings of the Association for Computational Linguistics: NAACL 2024, pp. 547 3846-3868, 2024. 548 Hila Gonen, Srini Iyer, Terra Blevins, Noah A Smith, and Luke Zettlemoyer. Demystifying prompts 549 in language models via perplexity estimation. In Findings of the Association for Computational 550 Linguistics: EMNLP 2023, pp. 10136–10148, 2023. 551 552 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented 553 language model pre-training. In International conference on machine learning, pp. 3929–3938. 554 PMLR, 2020. 555 Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop 556 ga dataset for comprehensive evaluation of reasoning steps. In Proceedings of the 28th Interna-557 tional Conference on Computational Linguistics, pp. 6609–6625, 2020. 558 559 Tenghao Huang, Ehsan Qasemi, Bangzheng Li, He Wang, Faeze Brahman, Muhao Chen, and Snigdha Chaturvedi. Affective and dynamic beam search for story generation. In *Findings of* 560 the Association for Computational Linguistics: EMNLP 2023, pp. 11792–11806, 2023. 561 562 Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open 563 domain question answering. In Proceedings of the 16th Conference of the European Chapter 564 of the Association for Computational Linguistics: Main Volume, pp. 874-880. Association for 565 Computational Linguistics, 2021. 566 Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand 567 Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning. 568 arXiv preprint arXiv:2112.09118, 2021. 569 Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili 570 Qiu. Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt com-571 pression. arXiv preprint arXiv:2310.06839, 2023a. 572 573 Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, 574 Jamie Callan, and Graham Neubig. Active retrieval augmented generation. In Proceedings of the 575 2023 Conference on Empirical Methods in Natural Language Processing, pp. 7969–7992, 2023b. 576 Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaga: A large scale distantly 577 supervised challenge dataset for reading comprehension. In Proceedings of the 55th Annual Meet-578 ing of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 1601–1611, 579 2017. 580 Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Dangi 581 Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In 582 Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing 583 (EMNLP), pp. 6769–6781, 2020. 584 585 Zixuan Ke, Weize Kong, Cheng Li, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. Bridg-586 ing the preference gap between retrievers and llms. arXiv preprint arXiv:2401.06954, 2024. 587 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris 588 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: A 589 benchmark for question answering research. Transactions of the Association for Computational 590 Linguistics, 7:452-466, 2019. 591 Yoonsang Lee, Pranav Atreya, Xi Ye, and Eunsol Choi. Crafting in-context examples according 592 to lms' parametric knowledge. In Findings of the Association for Computational Linguistics: 593 NAACL 2024, pp. 2069–2085, 2024.

594 595 596 597	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera- tion for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33: 9459–9474, 2020.
598 599 600 601	Bangzheng Li, Ben Zhou, Xingyu Fu, Fei Wang, Dan Roth, and Muhao Chen. Famicom: Further demystifying prompts for language models with task-agnostic performance estimation. <i>arXiv</i> preprint arXiv:2406.11243, 2024a.
602 603 604 605	Bangzheng Li, Ben Zhou, Fei Wang, Xingyu Fu, Dan Roth, and Muhao Chen. Deceptive semantic shortcuts on reasoning chains: How far can models go without hallucination? In <i>Proceedings</i> of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 7668–7681, 2024b.
606 607 608 609	Miaoran Li, Baolin Peng, Michel Galley, Jianfeng Gao, and Zhu Zhang. Self-checker: Plug-and- play modules for fact-checking with large language models. In <i>Findings of the Association for</i> <i>Computational Linguistics: NAACL 2024</i> , pp. 163–181, 2024c.
610 611 612 613	Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori B Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 12286–12312, 2023.
614 615 616 617	Zhonghao Li, Xuming Hu, Aiwei Liu, Kening Zheng, Sirui Huang, and Hui Xiong. Refiner: Re- structure retrieval content efficiently to advance question-answering capabilities. <i>arXiv preprint</i> <i>arXiv:2406.11357</i> , 2024d.
618 619 620	Qin Liu, Fei Wang, Nan Xu, Tianyi Yan, Tao Meng, and Muhao Chen. Monotonic paraphrasing improves generalization of language model prompting. <i>arXiv preprint arXiv:2403.16038</i> , 2024.
621 622 623	Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. Entity-based knowledge conflicts in question answering. In <i>Proceedings of the 2021 Conference</i> <i>on Empirical Methods in Natural Language Processing</i> , 2021.
624 625 626	Keming Lu, I-Hung Hsu, Wenxuan Zhou, Mingyu Derek Ma, and Muhao Chen. Multi-hop evidence retrieval for cross-document relation extraction. In <i>The 61st Annual Meeting Of The Association For Computational Linguistics</i> , 2023.
627 628 629 630	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational</i> <i>Linguistics (Volume 1: Long Papers)</i> , pp. 9802–9822, 2023.
632 633	Tao Meng, Sidi Lu, Nanyun Peng, and Kai-Wei Chang. Controllable text generation with neurally- decomposed oracle. <i>Advances in Neural Information Processing Systems</i> , 35:28125–28139, 2022.
634 635 636 637	Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. Document ranking with a pre- trained sequence-to-sequence model. In <i>Findings of the Association for Computational Linguis-</i> <i>tics: EMNLP 2020</i> , pp. 708–718, 2020.
638 639 640 641	Liangming Pan, Xiaobao Wu, Xinyuan Lu, Anh Tuan Luu, William Yang Wang, Min-Yen Kan, and Preslav Nakov. Fact-checking complex claims with program-guided reasoning. In <i>Proceedings</i> of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 6981–7004, 2023.
642 643 644	Nils Reimers and Iryna Gurevych. Making monolingual sentence embeddings multilingual using knowledge distillation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> . Association for Computational Linguistics, 2020.
646 647	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. Large language models can be easily distracted by irrelevant context. In <i>International Conference on Machine Learning</i> , pp. 31210–31227. PMLR, 2023.

660

668

684

648	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemover, and Wen-tau Yih.
649	Trusting your evidence: Hallucinate less with context-aware decoding. In <i>Proceedings of the 2024</i>
650	Conference of the North American Chapter of the Association for Computational Linguistics:
651	Human Language Technologies (Volume 2: Short Papers), pp. 783–791, 2024.
652	

- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Musique: Multihop
   questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554, 2022.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Interleaving re trieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 10014–10037, 2023.
- Hitesh Wadhwa, Rahul Seetharaman, Somyaa Aggarwal, Reshmi Ghosh, Samyadeep Basu,
  Soundararajan Srinivasan, Wenlong Zhao, Shreyas Chaudhari, and Ehsan Aghazadeh. From rags
  to rich parameters: Probing how language models utilize external knowledge over parametric
  information for factual queries. *arXiv preprint arXiv:2406.12824*, 2024.
- Fei Wang, Wenjie Mo, Yiwei Wang, Wenxuan Zhou, and Muhao Chen. A causal view of entity bias in (large) language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 15173–15184, 2023a.
- Yu Wang, Nedim Lipka, Ryan A Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. Knowledge graph
   prompting for multi-document question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 19206–19214, 2024.
- Zezhong Wang, Luyao Ye, Hongru Wang, Wai Chung Kwan, David Ho, and Kam-Fai Wong. Read-prompt: A readable prompting method for reliable knowledge probing. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 7468–7479, 2023b.
- Zhiruo Wang, Jun Araki, Zhengbao Jiang, Md Rizwan Parvez, and Graham Neubig. Learning to
   filter context for retrieval-augmented generation. *arXiv preprint arXiv:2311.08377*, 2023c.
- Kevin Wu, Eric Wu, and James Zou. Clasheval: Quantifying the tug-of-war between an llm's internal prior and external evidence. *Preprint*, 2024.
- Fangyuan Xu, Weijia Shi, and Eunsol Choi. Recomp: Improving retrieval-augmented lms with
   context compression and selective augmentation. In *The Twelfth International Conference on Learning Representations*, 2024.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2369–2380, 2018.
- Chanwoong Yoon, Taewhoo Lee, Hyeon Hwang, Minbyul Jeong, and Jaewoo Kang. Compact: Compressing retrieved documents actively for question answering. *arXiv preprint* arXiv:2407.09014, 2024.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, S Sanyal, Chenguang Zhu,
   Michael Zeng, and Meng Jiang. Generate rather than retrieve: Large language models are strong
   context generators. In *International Conference on Learning Representations*, 2023.
- Hao Zhang, Yuyang Zhang, Xiaoguang Li, Wenxuan Shi, Haonan Xu, Huanshuo Liu, Yasheng
  Wang, Lifeng Shang, Qun Liu, Yong Liu, et al. Evaluating the external and parametric knowledge
  fusion of large language models. *arXiv preprint arXiv:2405.19010*, 2024.
- Yunxiang Zhang, Muhammad Khalifa, Lajanugen Logeswaran, Moontae Lee, Honglak Lee, and Lu Wang. Merging generated and retrieved knowledge for open-domain qa. In *Proceedings of the* 2023 Conference on Empirical Methods in Natural Language Processing, pp. 4710–4728, 2023.

- Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and Muhao Chen. Context-faithful prompting for large language models. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, pp. 14544–14556, 2023.
- Honglei Zhuang, Zhen Qin, Rolf Jagerman, Kai Hui, Ji Ma, Jing Lu, Jianmo Ni, Xuanhui Wang, and
   Michael Bendersky. Rankt5: Fine-tuning t5 for text ranking with ranking losses. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 2308–2313, 2023.
- 710 711

712

702

703

704

705

#### A IMPLEMENTATION DETAILS

713 (1) Gold Compression: We implement the Gold Compression baseline following the approach 714 outlined by Yoon et al. (2024). We evaluate only on HotpotQA, 2WikiMQA, and MuSiQue, as 715 these datasets contain gold documents. We first identify the presence of any gold documents in 716 the retrieved documents. If found, we use the documents as the context. If none of the retrieved 717 documents are identified as gold, we utilize the entire set of retrieved documents as the context for 718 the evaluation. To identify the gold documents within the retrieved documents, we compare each 719 gold document with the retrieved ones. If 50% or more of the content matches, we classify it as a 720 gold document. This approach is necessary because the documents are chunked, and the retrieved documents may not exactly match the gold documents. 721

(2) **Generated Context**: We use the context generation prompt in Tab. 6 to generate the context.

723
 724 (3) Zero-shot Summarization: We use the evidence compression prompt in Tab. 6 to compress the retrieved documents.

(4) **RECOMP-extractive**: We utilize the same Contriever models trained by the authors for each dataset, to encode both the question and the sentences in the evidence set. For 2WikiMQA and MuSiQue, since there are no fine-tuned models available, we use the Contriever fine-tuned on HotpotQA. Following the original paper, we select one sentence as the context for NQ and TQA, whereas for the other datasets, we utilize two sentences.

(5) **RECOMP-abstractive**: Similar to RECOMP-extractive, we use the same T5-large models
trained by the authors for each dataset to compress the retrieved evidence. For the 2WikiMQA
and MuSiQue, we employ the T5-large model fine-tuned on HotpotQA.

(6) LongLLMLingua: We use Llama2-7B<sup>5</sup> trained by the authors as the prompt compressor model. We use the default hyperparameters in the original paper, where the dynamic context compression rate is set to 0.3, and the maximum compression rate is set to 0.5.

(7) CompAct: We use the same Mistral-7B-Instruct<sup>6</sup> model instruction-tuned by the authors for evidence compression. The number of documents per segment is set to 5 with 1 iteration.

740 741

742

755

#### **B** ADDITIONAL EXPERIMENT RESULTS

743 B.1 MISTRAL-7B-INSTRUCT AS COMPRESSION AND TARGET MODEL

We conduct an experiment where we use Mistral-7B-Instruct as the compression and target
 model. The result in Tab. 3 demonstrates that FAVICOMP outperforms all other baselines, supplementing the effectiveness shown in §4.1

748 749 B.2 HEAD-TO-HEAD COMPARISON WITH RECOMP-ABSTRACTIVE

We conduct a head-to-head experiment on RECOMP-abstractive by using the same base compression model as FAVICOMP for a more fair comparison. We construct training data on NQ, TQA, and HotpptQA according to Xu et al. (2024) and finetune Mistral-7B-Instruct on each of the training data. We train for 7 epochs using LoRA with Adam optimizer with a learning rate of

<sup>&</sup>lt;sup>5</sup>https://huggingface.co/NousResearch/Llama-2-7b-hf

<sup>&</sup>lt;sup>6</sup>https://huggingface.co/cwyoon99/CompAct-7b

Methods	Size	N	Q	T	QA	Hotp	otQA	2Wik	iMQA	MuS	SiQue
i i cuivas	Size	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
		Mi	stral-71	3-Instru	ct						
Gold Document		-	-	-	-	41.0	50.5	38.1	40.3	9.6	15.2
No Context Raw Document		28.1	27.5 39.3	58.8 66.2	60.9 68.6	19.7 30.3	24.8 37.2	21.9	22.8 28.5	5.2	9.7 13.1
Generated Context		30.1	31.7	57.3	60.7	23.7	30.6	25.1	29.5	7.1	12.
Sentence-BERT	110M	29.8	30.1	57.8	60.7	23.8	30.3	22.9	24.7	7.5	12.3
RECOMP-extractive	$110M^{\dagger}$	31.7	32.2	57.2	60.0	24.1	30.2	23.2	24.4	7.4	12.
LongLLMLingua	$7B^{\dagger}$	34.3	36.4	63.8	66.9	27.0	34.7	25.5	28.0	7.1	13.0
RECOMP-abstractive	775M <sup>†</sup>	38.0	37.8	62.1	65.0	27.4	34.3	25.1	27.4	6.4	12.0
CompAct	$7B^{\dagger}$	38.8	38.9	65.1	67.1	30.2	37.1	24.9	27.6	8.2	13.0
Zero-shot Summarization	7B	38.4	38.2	62.3	64.8	28.2	35.2	23.2	27.1	6.8	11.
FAVICOMP	7B	40.3	40.4	65.9	68.9	32.0	40.5	29.7	35.1	9.2	15.

Table 3: Experimental results when FAVICOMP has different compression and target models. We test
using Mixtral-8x7B-Instruct as the target model on five open-domain QA datasets across
all the methods. Mistral-7B-Instruct is used as the compression model of FAVICOMP. The
best Accuracy and token-level F1 scores for each dataset are in bold.

2e-6 and a batch size of 64. We present the evaluation results in Tab. 4. Even though using larger base model for compression enhances the performance of RECOMP-abstractive to some extent, it still underperforms compared to training-free FAVICOMP. This underscores that the familiarization during evidence compression and integration of parametric and non-parametric knowledge are more helpful to the downstream generation than relying on a trained model for evidence compression.

Methods	Train	Compression Model	N	Q	T(	)A	Hotp	otQA
		<b>F</b>	Acc	F1	Acc	F1	Acc	F1
RECOMP-abstractive	0	T5-large	38.0	37.8	62.1	65.0	27.4	34.3
<b>RECOMP-abstractive</b>	0	Mistral-7B-Instruct-v0.3	38.3	38.2	63.0	65.4	29.5	36.6
FaviComp	Х	Mistral-7B-Instruct-v0.3	40.3	40.4	65.9	68.9	32.0	40.5

Table 4:	Head-to-head	comparison	results	with RECOMP
14010	rieda to nead	• ompanioon	1000100	n null rub e e null

#### C PROMPT TEMPLATES

#### C.1 EVALUATION

The evaluation prompt template is shown in Fig. 5. For all the evaluations throughout the experiment, we switch the positions of the Question and Context if doing so results in better performance.
System prompts and demonstrations used in the evaluation are presented in Tab. 5 and Tab. 7, respectively.

#### C.2 FAVICOMP

The prompt templates for evidence compression and context generation of FAVICOMP are presented in Tab. 6





**Evaluation Prompt Template** {System Prompt} {Demonstrations} Question: {Question} Context: {Context} Answer:

#### Figure 5: Evaluation Prompt Template.

Model	System Prompt
Llama-3-8B-Instruct	You are an expert in Question Answering. Your job is to answer questions in 1 to 5 words based on the given context.
Mistral-7B-Instruct	You are an expert in Question Answering. Your job is to answer questions in 1 to 5 words based on the given context. Just output the answer as concisely as possible, no other words
	Table 5: System prompts used in evaluation
	Table 5. 5ystem prompts used in evaluation

865		
866		
867		
868		
869		
870		
871		
872		
873		
874		
875		
876		
877		
878		
879		
880		
881		
882		
883		
884		
885	Instruction	Prompt Template
886		Vou are an expert in summerization. Given a question and multiple decument enimeter
887		generate one summarized context that is helpful to answer the question. Just summa-
888	Evidence Commercian	rize, no other words.
889	Evidence Compression	Question: {Question}
890		Documents: {Evidence}
891		Summarized Context:
892		You are an expert in context generation. Given a question, generate a context that is
893	Context Generation	helpful to answer the question. Just generate the context, no other words.
894		Context:
895		I
896		Table 6: Prompt Templates for FAVICOMP
897		
898		
000		
900		
902		
903		
904		
905		
906		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		

Dataset	Demonstrations
	Question: who sings i've got to be me
	Answer: Sammy Davis, Jr
	Question: who wrote i will follow you into the dark
	Answer: Ben Gibbard
VQ	Answer: Owen (Scott McCord)
	Question: what part of the mammary gland produces milk
	Answer: cuboidal cells
	Question: when did the golden compass book come out
	Our diam What some the drawe for the Lance Dend film (Thur dashell')
	Answer: Tom Jones
	Question: A hendecagon has how many sides?
	Answer: Eleven
	Question: In the 1968 feature film Chitty Chitty Bang Bang, of what country is Baron Bomburst the tyrant ruler?
TQA	Answer: Vulgaria
	Question: Artists Chuck Close, Henri-Edmond Cross, John Roy, Georges-Pierre Seurat, Pau
	Signac, Maximilien Luce and Vincent van Gogh painted in what style?
	Ouestion: What is the study of the relation between the motion of a body and the forces acting of
	it?
	Answer: Dynamics
	Question: Which magazine was started first Arthur's Magazine or First for Women?
	Answer: Arthur's Magazine
	Answer: Delhi
	Question: Musician and satirist Allie Goertz wrote a song about the "The Simpsons" characte
HotpotQA	Milhouse, who Matt Groening named after who?
	Answer: President Richard Nixon Question: Are Jane and First for Women both women's magazines?
	Answer: Yes
	Question: Were Pavel Urysohn and Leonid Levin known for the same type of work?
	Answer: No
	Question: Where was the place of death of Marie Thérèse Of France (1667–1672)'s father?
	Question: Who is the paternal grandmother of Przemysław Potocki?
	Answer: Ludwika Lubomirska
	Question: Who lived longer, Herbert Findeisen or Léonie Humbert-Vignot?
2 WIKIMQA	Ouestion: Are Alison Skipper and Diane Gilliam Fisher from the same country?
	Answer: Yes
	Question: Are director of film Move (1970 Film) and director of film Méditerranée (1963 Film
	from the same country?
	Question: who is the child of the director and star of Awwal Number? Answer: Supeil Anand
	Question: What is the record label of the rapper who performed Jigga My?
	Answer: Roc-A-Fella Records
	Question: What county shares a border with the county where Black Hawk Township is located?
MuSiQue	Ouestion: Who is the sibling of the person credited with the reinvention and popularization of oi
	paints?
	Answer: Hubert Van Eyck
	Question: who heads the Catholic Church, in the country that a harp is associated with, as a lion is associated with the country that Queen Margaret and her son traveled to?
	Answer: Eamon Martin
	1
	Table 7: Demonstrations used in evaluation for each dataset