# INVAR-RAG: INVARIANT LLM-ALIGNED RETRIEVAL FOR BETTER GENERATION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Retrieval-augmented generation (RAG) has shown its impressive capability of providing reliable answer predictions and addressing severe hallucination problems. A typical RAG implementation adopts powerful retrieval models to extract external information and leverage large language models (LLMs) to generate corresponding answers. Different with that, recent LLM-based retrieval has raised much attention because it brings substantial improvements in information retrieval (IR) via LLMs' vigorous semantic understanding capability. However, directly applying LLM to RAG systems remains certain challenges. This may cause feature locality problems since massive parametric knowledge impedes the effective usage of the global information among all corpus, e.g., a LLM-based retriever usually inputs the summary of documents instead of the whole documents. Moreover, various tasks pre-trained in LLMs induce severe variance, which further weakens its performance as the retriever. To address these issues, we propose a novel two-stage fine-tuning architecture called Invar-RAG. In the retrieval stage, a LLM-based retriever is constructed by integrating a LoRA-based representation learning to address the feature locality problem. To justify and consolidate this retrieval's performance, two patterns (i.e., invariant and variant patterns) and an invariance loss are also developed to alleviate the variance in LLM. Moreover, in the generation stage, a meticulously designed fine-tuning method is devised to improve our LLM for accurate answer generation based on the retrieved information. Experimental results demonstrate that Invar-RAG significantly outperforms existing baselines across three Open-domain Question Answering (ODQA) datasets. The code is available in Supplementary Material to ease reproducibility.

032 033 034

035

043

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

#### 1 INTRODUCTION

Over the past decade, large language models (LLMs) have demonstrated promising capability in processing natural language Minaee et al. (2024). Owing to the vast amount of knowledge encoded in their internal parameters, LLMs such as GPT Achiam et al. (2023) and LLaMa Touvron et al. (2023) have demonstrated remarkable performance on various downstream tasks, including Open-domain Question Answering (ODQA) Zhu et al. (2021), Reading Comprehension Cheng et al. (2023). However, the fixed parametric knowledge of LLMs has hindered the further applications of LLMs and made them prone to errors (hallucination Rawte et al. (2023) and factual errors Wang et al. (2023)).

To overcome the limitations of parametric knowledge, one promising approach is Retrieval-044 Augmented Generation (RAG) Wang et al. (2023); Lewis et al. (2020). Compared to relying solely on parametric knowledge, RAG enables LLMs to use retrievers to access relevant information from 046 external knowledge sources, enhancing their question-answering abilities. Among the two compo-047 nents of RAG, current methods primarily focus on optimizing the retriever to return more relevant 048 documents due to the high cost of fine-tuning and black-box LLM APIs. Previous retrievers leveraged deep learning technology (e.g., dense retrieval Zhao et al. (2024)) to encode the text representations from the lexical space into the high-dimensional latent space, allowing them to model more 051 complex semantic relationships between queries and corpora. However, the separation between the retriever and generation model has hindered their full integration, limiting their compatibility in 052 downstream applications. Some advanced RAG systems, such as RA-DIT Lin et al. (2023b), have adopted joint training mechanisms that fine-tune both the retriever and the generation model for better alignment. However, this approach is impractical due to the need for frequent fine-tuning and
 fails to utilize the LLMs' semantic understanding capabilities during the retrieval stage.

Consequently, generative retrieval (GR), also known as LLM-based retrieval, leverages the para-057 metric memory of generative models to directly generate document identifiers (DocIDs) Li et al. (2024), which has aroused much attention. By memorizing the documents as the parametric knowledge of LLM, this kind of method breaks the limitations of traditional IR in terms of document 060 granularity and simple relevance matching Nguyen & Yates (2023), offering more flexibility and 061 creativity, thus better meeting practical needs. However, two severe problems hinder the current 062 LLM-based retrieval. 1)Feature Locality: LLM-based retrieval normally adopt language models 063 to learn the mapping from queries to the relevant document DocIDs. However, these DocIDs ac-064 tually can not fully represent the global information of the passages. Meanwhile, directly feeding the whole passages into LLM is costly and infeasible, causing a trade-off between effectiveness and 065 efficiency. 2)Retrieval Variance: Due to the inherent generative inconsistency property of large 066 language models, current LLM-based retrieval may generate unforeseeable variances, especially 067 when the input query or the size of context varies, directly results in undesirable and vulnerable 068 performance which may not be preferred. 069

Considering the problems mentioned above and better leveraging the capability of LLMs, we pro-071 pose a fully LLM-based architecture with a two-stage fine-tuning method called Invar-RAG, as illustrated in Figure 1. In the retrieval stage, our approach initializes the pre-trained LLaMA Tou-072 vron et al. (2023) as the backbone and follows the bi-encoder architecture in DPR Karpukhin et al. 073 (2020) to construct our retriever. Compared to normal GR methods which need an iterative pro-074 cess of encoding and decoding, we introduce a component called LLM-aligned Retrieval. It first 075 represents the input query and corpora into high-dimension space using a small language model 076 (MiniLM) Wang et al. (2020), then introduces a new loss function constructed by KL-divergence 077 to align the coarse query-documents pairs representation to the LLM's representation space. This allows the retriever to leverage the rich prior knowledge of LLM, typically addressing the feature 079 locality caused by only feeding DocIDs to LLM. Moreover, based on the initial objective of our LLM-aligned Retrieval, we introduce the invariance loss to overcome the variance in the retrieval 081 stage. By recognizing the invariant pattern that contributes the most to the performance and gradually forcing the model to rely on the invariant pattern, we can avoid the unforeseeable variances in 082 practice and enhance the robustness of our RAG system. Finally, in the generation stage, we freeze 083 the weights we fine-tuned before and optimize the generation function to allow the LLM to give 084 correct answers to the retrieved documents. Our contributions are summarized as follows: 085

- We introduce Invar-RAG, a novel framework featuring a two-stage fine-tuning method on a single shared LLM, including the retrieval stage and generation stage.
- We introduce a novel LLM-based retrieval method containing representation learning and invariance loss, respectively addressing the issues of feature locality and retrieval variance.
- We validate Invar-RAG's performance on three public ODQA datasets, no matter for retrieval performance or generation performance, demonstrating its superiority.

## 2 Methodology

In this section, we introduce a novel retrieval-augmented language model architecture, Invar-RAG, which addresses the previously mentioned issues by using LLM-aligned retrieval combined with a specially designed invariance loss. We first present an overview of our proposed architecture, followed by a detailed explanation of its key components, and finally, we introduce how we construct the prompts.

101 102 103

090

092 093 094

095 096

097

098

099

100

#### 2.1 OVERALL FRAMEWORK OF INVAR-RAG

In this section, we provide an overview of Invar-RAG, as shown in Fig. 1. We begin by using query rewriting and context window resizing to introduce various types of variance. Next, we apply a small language model to map these texts into a high-dimensional vector space, generating coarse representations. We then adopt LLM-aligned retrieval to align the coarse representation with the LLM's representation and compute the basic relevance score via dot product. This approach



Figure 1: Overview of proposed Invar-RAG.

addresses the feature locality problem by feeding the entire document representation into the LLM, 129 rather than a single DocID. Additionally, to address feature variance, we define an invariance loss based on the initial KL-divergence loss function in representation learning, encouraging the model to rely on invariant patterns. Finally, by constructing appropriate prompts and fine-tuning the generation stage, we optimize the LLM to better utilize the retrieved information, generating more accurate answers to the given questions.

137

142

126

127 128

130

131

132

#### 2.2 **RETRIEVAL STAGE**

Architecture. For the Retriever architecture, we follow the approach of previous work Ma et al. 138 (2024), using the bi-encoder architecture from DPR Karpukhin et al. (2020), but replacing the back-139 bone model with LLaMA Touvron et al. (2023). Considering the efficiency, we first compute the 140 vector embedding of a document  $d_i^R \in \mathcal{D}^{\hat{R}}$  as: 141

$$V_r(d_i^R) = \text{Decoder}\left({}^t t_1 \lrcorner t_2 \lrcorner \cdots \lrcorner t_k'\right) [-1] \tag{1}$$

143 Where k represents the maximum number of trunks, and  $Decoder(\cdot)$  represents the embedding layer 144 of a small language model (MiniLM-v2), which maps the trunks  $(t_1 \downarrow t_2 \sqcup \cdots \downarrow t_k)$  from the initial text 145 space to a high-dimensional dense vector space. For the vector embedding of query q, we leverage 146 our LLM structure to return the last layer token representation as the representation, denoted as  $V_q$ . 147

To leverage the LLM's prior knowledge while maintaining efficiency, we further align the query-148 document pairs to the LLM's representation space and denote the processed document representation 149 as  $V_p(d_i^R)$ . Consequently, we can compute in terms of the dot product to get the two relevance 150 scores: 151

152 153

 $s_{raw}\left(V_q, V_r(d_i^R)\right) = V_q \cdot V_r(d_i^R)$ (2) $s_{pro}\left(V_q, V_p(d_i^R)\right) = V_q \cdot V_p(d_i^R)$ 

where the basic relevance score between the query and documents processed by small LM denotes 154 as  $P_{raw}$  and the target relevance score computing between query and LLM-processed documents 155 denotes as  $P_{pro}$ . 156

157 LLM-aligned Retrieval. Building on the above structure, we enhance our LLM-based retriever's 158 ability to return more relevant documents. We introduce LLM-aligned retrieval with invariance loss 159 in the retrieval stage, effectively addressing the aforementioned issues. Current alignment methods, such as RA-DIT Lin et al. (2023b), primarily focus on aligning the scoring functions between the 160 retriever and the generator. However, the initial structural differences between novel retrievers (e.g., 161 DRAGON+ Lin et al. (2023a)) and LLMs still impede further optimization of the overall RAG 169 170

171 172

177

193 194

196 197

199

207

208

212 213

system. Therefore, we design a novel LLM-based retriever to resolve this issue. Unlike previous LLM-based GR methods Li et al. (2024), we no longer need to use DocIDs to retrieve relevant documents, which may cause the feature locality issue mentioned earlier. Instead, we adopt a finetuned alignment process that enables the LLM to perform representation learning. We leverage LoRA architecture Hu et al. (2021) to add additional adapter parameter  $\theta_R$  to our raw representation  $V_r$ , denoted as  $V_r(d_i^R, \theta_R)$ . The corresponding relevance score can then be re-normalized among top-k relevant chunks  $\mathcal{D}^{R'} \subset \mathcal{D}^R$  as:

$$S_{\mathcal{R}}^{r}(V_{r}(d_{i}^{R},\theta_{\mathcal{R}})|V_{q}) = \frac{\exp s_{raw}\left(V_{q},V_{r}(d_{i}^{R},\theta_{\mathcal{R}})\right)}{\sum_{d_{i}^{R'}\in\mathcal{D}^{R'}}\exp s_{raw}\left(V_{q},V_{r}(d_{i}^{R'},\theta_{\mathcal{R}})\right)}$$
(3)

For each document in the corpus, we need to compute  $S_{\mathcal{R}}^r$  for *n* times (*n* represents the number of documents in  $\mathcal{D}^R$ ) to rank the relevance scores. Consequently, the initial loss function for representation learning can then be defined by minimizing the KL-divergence Kim et al. (2021) of two relevance scores leveraging Eq. 3.

$$\mathcal{L}_{rl}\left(\mathcal{D}^{R}\right) = \mathbb{E}_{d_{i}^{R}\in\mathcal{D}^{R}}\operatorname{KL}\left(S_{\mathcal{R}}^{r}(V_{r}(d_{i}^{R},\theta_{\mathcal{R}})|V_{q})\|S_{\mathcal{R}}^{r}(V_{p}(d_{i}^{R})|V_{q})\right)$$
(4)

Following the previous works Lin et al. (2023b); Ma et al. (2020), fine-tuning both encoders hurt the performance Bao et al. (2024), we only update a part of our initialized retriever, which is in charge of computing the query representation.

181 Invariance Loss. To further enhance retrieval accuracy while maintaining robustness, we introduce 182 invariance loss, building on our initial KL-divergence loss. Current refinement methods primarily 183 rely on query rewriting He et al. (2016); Chan et al. (2024) or LLM generation Fan et al. (2024) to 184 expand the search space and re-rank document chunks. However, they fail to recognize the effective-185 ness of different rewriting and generation procedures, directly resulting in the invariance problem. Specifically, we begin by rewriting the query and adjusting the context window to broaden the search 186 space. Since not all refinement methods are effective, we identify invariant patterns to preserve re-187 trieval performance while gradually incorporating weighted variant patterns to broaden the search 188 space. To achieve this, we use the LSR score from LM-Supervised Retrieval Shi et al. (2023) to 189 determine whether a document effectively enhances the LLM's answer prediction capability. For a 190 training sample (q, y), where q and y respectively represent the input query and output result, we 191 first define the output probability of LM as: 192

$$p_{LM}(y|V(d_i^R \circ x)) = \sum_{d_i^R \in \mathcal{D}^R} p_{LM}(y|V(d_i^R \circ x)) \cdot P_{\mathcal{R}}(d_i^R|x)$$
(5)

Then, for the LSR score for a retrieved document  $d_i^R$ :

$$P_{LSR}(d_i|q,y) = \frac{\exp(p_{LM}(y|d_i \circ q)/\tau)}{\sum_{d'_i \in \mathcal{D}^{R'}} \exp(p_{LM}(y|d'_i \circ q)/\tau)} \approx \frac{\exp(p_{LM}(y|d_i \circ q)/\tau)}{\sum_{d'_i \in \mathcal{D}^R} \exp(p_{LM}(y|d'_i \circ q)/\tau)} \quad (6)$$

where  $\tau$  is the temperature hyperparameter of LLM,  $\mathcal{D}'_R \subset \mathcal{D}_R$  denotes the top-k retrieved trunks. Assuming the query after rewriting as  $q^r$ , documents set after resizing as  $\mathcal{D}_R^{re} = (d_1^{re}, d_2^{re}, \cdots, d_n^{re})$ , we can leverage the Eq. 6 to calculate the score matching from (1) q to  $d_i$ , (2)  $q^r$  to  $d_i$ , (3) q to  $d_i^{re}$ and (4)  $q^r$  to  $d_i^{re}$ . We recognize the invariant pattern as the top-l ranked documents, denoted as  $\mathcal{D}_{in}$ , where 0 < l < k is satisfied. For other documents, we assume them as variant pattern  $\mathcal{D}_{var}$ , which contribute little to generating effective answers.

206 The invariant loss function can be formalized as follows:

$$\mathcal{L}_{invar}(\mathcal{D}_{in}) = \operatorname{Var}_{\mathcal{D}\subseteq\mathcal{D}_{var}}(\mathbb{E}_{d_{in}^R\in(\mathcal{D}_{in}\cup\mathcal{D})}\operatorname{KL}(P_{\mathcal{R}}^r(V_r(d_{in}^R,\theta_{\mathcal{R}})|V_q)\|P_{\mathcal{R}}^p(V_p(d_{in}^R)|V_q)))$$
(7)

This invariance loss measures the variance of the model's aligning ability under multiple interventions (*i.e.*, query rewriting and context resizing) by only allowing the documents in the  $\mathcal{D}_{invar}$  to update the loss. The whole training objective can then be presented as:

$$\min_{\alpha} \mathcal{L}_{rl} + \lambda \mathcal{L}_{invar} \tag{8}$$

where the task loss  $\mathcal{L}_{rl}$  is minimized to align the two different representations while the  $\mathcal{L}_{invar}$ enables the model to rely more on the invariant pattern, and  $\lambda$  is a hyperparameter to balance between two objectives.

4

216	Table 1: The Statistics of Fine-tuning Datasets.								
217	Dataset	HF identifier	$\mathcal{D}_R$	$\mathcal{D}_G$	Training Sample	Task			
218	Wiki QA Yang et al. (2015)	wiki_qa	X	1	20360	Open-domain QA			
219	FreebaseQA Yao et al. (2014)	freebase_qa	1	X	20358	Open-domain QA			
220	MS-MARCO Bajaj et al. (2016)	ms_marco	1	X	80143	Open-domain QA			
221	Web Question Dumais et al. (2002)	web_question	X	1	3778	Open-domain QA			
222	SQuAD v2	squad_v2	X	1	130319	Reading Comprehension			

Table 1: The Statistics of Fine-tuning Datasets.

#### 2.3 GENERATION STAGE

224

235

236

242 243

244

249

251

252

260

262

263

264

265 266

267

268

269

225 In the generation stage, We followed the same architecture as in the retrieval stage for answer predic-226 tion. To improve the generative capability of LLM for leveraging the retrieved information better, 227 followed by prior works Lin et al. (2023b); Shi et al. (2023), we adopt another LoRA adapter to 228 fine-tune our model on different tasks. Specifically, for the same training sample (x, y), we retrieve 229 the *top-k* relevant document chunks  $\mathcal{D}'_G \subset \mathcal{D}_G$  by performing our model on retrieval task. For 230 each retrieved chunk  $d_i \in \mathcal{D}'_G$ , we design a special fine-tuning example by prepending it to the 231 prompt as background information and create k independent instances for one original example: 232  $\{(d_i \circ x, y) | i = 1, \dots, \overline{k}\}$ . Then, following the previous work Lin et al. (2023b); Qi et al. (2020), 233 we fine-tune the language model using the next-token prediction objective and minimize the loss as 234 follows:

$$\mathcal{L}(\mathcal{D}'_G) = -\sum_i \log p_{LM}(y|d_i \circ x) \tag{9}$$

By applying this fine-tuning method, the generation stage benefits in two ways: (i) it improves
the model's performance on the generation task by providing more accurate predictions based on
the retrieved information; (ii) when the retrieved documents fail to provide an accurate answer,
the approach enables the LLM to rely on its parametric knowledge to generate an answer while
disregarding misleading retrieved documents.

#### 3 EXPERIMENT

In this section, we will first introduce the experiment setting. Then we present extensive experiments
to evaluate the effectiveness of our proposed Invar-RAG architecture in different stages (retrieval
and generation). All the reported experimental results are the average values obtained from five
independent runs of the algorithm.

250 3.1 SETTING

#### 3.1.1 DATASETS

Following the prior works Lin et al. (2023b); Asai et al. (2023), we choose two ODQA datasets (FreebaseQA Yao et al. (2014) and MS-MARCO Bajaj et al. (2016)) and one reading comprehension (RC) dataset to do the representation learning in the retrieval stage (denoted as  $D_R$ ) while leveraging other three datasets (Web Question Dumais et al. (2002), Wiki Question Answering Yang et al. (2015)) and SQuAD v2<sup>1</sup> for fine-tuning the LLM in the generation stage (denoted as  $D_G$ ). The statistic of chosen datasets is shown in Tab.1. For detailed descriptions and complied templates, please refer to **Appendix A**.

261 3.1.2 EVALUATION

To access our performance, we conduct the evaluation on four knowledge-intensive datasets, such as *i.e.*, TriviaQA (denoted as TQA)<sup>2</sup>, Natural Question (denoted as NQ)<sup>3</sup> and PopQA<sup>4</sup>, that are not involved in the training progress. For the evaluation metric, we evaluate our model's generation

<sup>1</sup>https://rajpurkar.github.io/SQuAD-explorer/explore/v2.0/dev/Prime\_ number.html

<sup>3</sup>https://ai.google.com/research/NaturalQuestions

<sup>4</sup>https://huggingface.co/datasets/akariasai/PopQA

<sup>&</sup>lt;sup>2</sup>https://nlp.cs.washington.edu/triviaqa

270 performance using the Exact Match Wongsuphasawat et al. (2012), which indicates whether gold an-271 swers are included in the model generations followed by the setting in prior work Lin et al. (2023b); 272 Mao et al. (2024). Furthermore, to evaluate our proposed retriever's performance, we employ the 273 Acc@5 and Acc@20 as evaluation metrics, which are widely used in related studies Chen et al. 274 (2024); Izacard et al. (2021). These metrics assess the proportion of questions where the correct answers appear in the top-5 or top-20 retrieval results, offering a comprehensive evaluation of the 275 retrieval performance. For more details, please refer to the description and methods in Appdendix 276 B. 277

278

#### 279 3.1.3 IMPLEMENTATION DETAILS

In this section, we provide a detailed description of our framework's implementation. The code can be found in Supplementary Material. For both the retrieval and generation stages, the LLaMA-2-7B checkpoint<sup>5</sup> is leveraged to initialize the pre-trained weights of our architecture. For the GPU selection, We perform our further fine-tuning on 4 × 40G NVIDIA V100 GPUs.

284 Retrieval Stage. Following the previous work's setting Ma et al. (2024), as LLaMA is a decoder-285 only architecture, we append an end-of-sequence token <EOS> to the input sequence and regard 286 the last layer representation as the dense representation to calculate the similarity score. Considering 287 the possible effect caused by the size of each dense representation, we also employ the normalization 288 procedure to map the original representation into unit vectors during both the training and inference 289 stages. For the fine-tuning progress in the retrieval stage, we adopt LoRA architecture Hu et al. 290 (2021) to reduce the high cost of GPU memory. The detailed hyperparameters we used can be found 291 in Appendix C.

Generation Stage. We hold  $\overline{k}$  in Sec.2.3 equal to 5 to generate instances for a single example and append multiple examples together to improve the efficiency (the length is limited to 4096 tokens). The used hyperparameters are also shown in the Appendix C. Other implementation details are the same as original papers Lin et al. (2023b); Shi et al. (2023).

# 297 3.1.4 BASELINES

To demonstrate the effectiveness of our proposed architecture, we compare the retrieval performance of our Invar-RAG with state-of-the-art retrieval methods, including sparse retrieval (BM25 Ram et al. (2023)), dense retrieval (BGE Xiao et al. (2024), Contriever Izacard et al. (2021)) and LLMbased retrieval (LLM-embedder Zhang et al. (2023) and RepLLaMA Ma et al. (2024)). Furthermore, for the corresponding RAG performance, we conduct extensive experiments compared to the novel retriever + generation model to show our superiority. The descriptions for each baseline are listed in the **Appendix D**.

305 306 307

#### 3.2 OVERALL PERFORMANCE

In this section, we present performance comparison experiments on two stages, respectively, with
 three knowledge-intensive ODQA datasets. The results show that our Invar-RAG architecture out performs all competing sparse, dense, and LLM-based baselines in retrieval and their downstream
 RAG in generation. Such a comparison highlights the effectiveness of our unique design for two stage fine-tuning.

313 314

#### 3.2.1 RETRIEVAL PERFORMANCE

315 In this section, we will present and analyze the retrieval performance of our designed architecture. 316 As illustrated in Tab.2, the sparse retriever BM25 fails to map the given text to proper representa-317 tions. Although employing an additional model as the re-ranker improves the performance to some 318 extent, the retrieval capability remains sub-optimal due to the inferiority of BM25. Besides, Novel 319 dense retrievers, like BGE and Contriever, present comparable performance over the three datasets, 320 suggesting their effectiveness in leveraging contrastive learning or task-specific fine-tuning. How-321 ever, they still slightly lag behind our designed Invar-retrieval because of the neglect of rich semantic 322 information Ma et al. (2024). Current researchers have proposed several LLM-based retrievers *i.e.*,

<sup>323</sup> 

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

325	Table 2: Retrieval performance comparison between our designed retriever in Invar-RAG and other
326	baselines. The best results are bold, and the second-best are underlined.

,	Madala	Т	QA	NQ		Poj	pQA
3	Models	Acc@5	Acc@20	Acc@5	Acc@20	Acc@5	Acc@20
)	BM25 Ram et al. (2023)	62.5	73.0	49.0	67.0	35.5	51.5
	BM25+BGE(re-ranker) Chen et al. (2024)	72.5	78.0	68.0	76.5	54.0	60.0
	Contriever Izacard et al. (2021)	68.0	80.5	68.0	84.0	62.0	77.5
	BGE-base Xiao et al. (2024)	69.5	80.0	77.0	86.0	72.0	83.0
	LLM-embedder Zhang et al. (2023)	67.5	77.5	75.5	86.5	70.0	79.5
	RepLLaMA Ma et al. (2024)	66.5	76.0	72.0	85.5	68.5	74.5
	Invar-retrieval (ours)	74.0	81.5	80.5	88.0	73.5	<u>82.5</u>
	Improv.	2.1%	1.2%	4.6%	1.7%	2.1%	-0.6%
i i							

Table 3: Generation performance comparison between our designed Invar-RAG and other baselines. The best results are bold, and the second-best are underlined.

339	Madala	TQA	PopQA	NQ
340	WIOdels	Exact Match		
341	BGE-base + LLaMA-2-7B	74.1	49.8	52.1
342	BM25 + BGE(re-rank) + LLaMA-2-7B	72.3	48.2	51.6
343	LLM-embeder + LLaMA-2-7B	71.8	51.1	54.1
344	Contriever + LLaMA-2-7B	72.6	48.6	51.8
345	Invar-RAG	75.3	53.6	56.2
346	Improv.	1.6%	4.9%	3.9%
347			, ,-	

348 LLM-embedder Zhang et al. (2023) and RepLLaMA Ma et al. (2024), which leverage the rich prior 349 knowledge that LLM initially has. However, due to the high cost of processing the massive corpus, it 350 is infeasible to handle all the chunks within the LLM. Moreover, the variance problem that happens 351 in LLM also leads to relatively inferior performance. Correspondingly, we propose our LLM-based 352 retrieval model, Invar-retrieval, as a part of our designed Invar-RAG. The results shows that our 353 methods outperform all the sparse, dense and LLM-based retrievers, especially under the Acc@5 measurement, contributing to our designed invariance loss in reducing the variant and ineffective 354 patterns. 355

356 357

324

337

338

#### 3.2.2 **GENERATION PERFORMANCE**

358 In this section, we will analyze the answer generation capability for our designed Invar-RAG. Based 359 on the astonishing performance of our designed Invar-retrieval, we further fine-tune the language 360 model to leverage the retrieved documents for better question-answering capability. From the ex-361 perimental results presented in Tab.3, our Invar-RAG shows reasonable performance on the three 362 ODQA datasets, echoing the performance of the retrievers designed in Tab.2.

363 364

#### 3.3 ABLATION STUDY

366 In this section, we analyze the efficacy of the two-stage fine-tuning in the Invar-RAG architecture, including the retrieval stage (LLM-aligned Retrieval with Invariance Loss) and the generation stage. 367 We design three variants: (1)w/o representation learning: this variant uses the coarse text representa-368 tion mapped by small language model (MiniLM-v2<sup>6</sup>) to calculate the relevance score and adopt the 369 same generation fine-tuning method in Sec.2.3. (2)w/o invariance loss: the second variant leverages 370 the KL-divergence loss without the additional invariance loss to perform the representation learning. 371 (3)w/o generative fine-tuning: this variant directly feeds retrieved documents and the corresponding 372 question as a prompt to generate the answer. The fine-tuning datasets for each variant we used are 373 presented in Tab.4. From the performance comparison in Tab.5, We can conclude that: 374

375 376

377

• With the representation learning method, LLM-based retrieval contributes to improving the retrieval and corresponding generation performance.

<sup>&</sup>lt;sup>6</sup>https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2

Model Variants	Retrieval Fin Freebase QA	e-tuning	O Wiki QA	Generation Fine-to Web Question	uning SQuAD v2
Default	✓	1	 ✓	· ·	/
w/o representation learn	ning 🗶	×	1	1	1
w/o invariance loss			1		<i>✓</i>
w/o generative fine-tuni	ing 🗸	1	X	X	X
	Table 5: Abl	ation Stud	y on TQA.		
١	Andal Varianta	Ret	rieval	Generation	
1	and variants	Acc@5	Acc@20	Exact Match	
Defaul	t	74.0	81.5	75.3	
w/o rep	presentation learning	63.5	73.6	74.1	
w/o inv	variance loss	71.5	81.0	74.6	
w/o ge	nerative fine-tuning	1	/	/3.4	
Invariance loss more on invari	s significantly boosts c ant patterns.	our designe	ed Invar-RA	G by making the	prediction rel
<ul> <li>Generative fin</li> </ul>	e-tuning is crucial fo	or enhancii	ng LLM's c	apability of givi	ing prediction
based on retrie	ved information. Mor	eover, it sh	nows the effe	ectiveness of the	two-stage fine
tuning for a sir	ngle LLM.				
7 11 . 1.			<b>X A</b> X <b>1</b>		P F
for ablation results on o	other two datasets (NC	2 and Pop	(A), please	found them in Ap	ppendix E.
0.4 INVARIANCE AN	ALYSIS				
n this section, we lever	age a special example	in TOA to	o illustrate tl	ne effectiveness o	of our designe
nvariance loss in two p	parts: (i) the importan	ce of defir	ning differer	t patterns, (ii) th	e difference i
etrieval performance th	at invariance loss brir	igs.			
As mentioned in Sec. 2	.2, we return four diff	ferent sets	of retrieved	documents and 1	rerank them b
LSR score to identify th	e invariant pattern. Th	here are tw	o reasons to	explain this:	
~					
• Rewriting the c	luery and resizing the	context wi	ndow does a	ffect the normal i	relevance scor
computing by	the dot product, lead	ng to the v $\mathbf{P} \mathbf{A} \mathbf{G}$ over	variance in j	or the answer	e feed differen
		e KAU sys		or the answer.	
Prior works Zi	ang et al. (2024) have	e shown th	at adding a	suitable amount	of irrelevant o
relatively men	ective documents does	s neip mip	love the let	ievai periormane	е.
To verify that, we prese	nt the normal relevand	ce score ar	nd LSR scor	e of each retrieve	ed document i
our different sets in Fi	g. 3.4. The darker co	olor represe	ents the cha	nge that happene	ed in the Top-
documents. We can se	e that, for the question	on: 'Who	was the ma	n behind The Cl	hipmunks', th
elevance score for the	top-5 documents in ea	ich set show	ws substanti	al changes while	the LSR scor
loes not vary a lot, which	ch means the variance	caused by	rewriting qu	uery or resizing c	ontext window
change the importance	of documents, directl	y resulting	g in poor rei	trieval performan	te our selecte
uestion example while	our designed Invar-R	AG system	accurately	predicts the result	lo our serecte lt
question example white	our designed nivar it	rio system	raccuratory	predicts the resu	
+ KELATED WOR	K				
		1. •	. 1	1 (1 1 . 1 . 6	
<b>Information Retrieval</b>	: Advancements in de	ep learnin	g have revo	iutionized inform	nation retrieva
mation retrieval framew	orks employed sparse	retrievers	Rametal (	2023) or dense ret	ns. Earry IIII01 trievers Izacar
et al. (2021): Xiao et al.	(2024) to represent lar	ge corpora	but struggle	ed to capture deer	semantic rela
tionships Karpukhin et a	al. (2020). LLM-based	l retrievers	(generative	retrieval) have sin	nce emerged a
otable methods, levera	ging the rich prior kno	wledge of	LLMs to sig	nificantly improv	ve performanc

Table 4. The Statistics of Datasets in Ablation study	Table 4:	The	Statistics	of <b>E</b>	Datasets	in	Ablation	study.
---	----------	-----	------------	-------------	----------	----	----------	--------

by converting documents into parametric knowledge and generating them instead of computing sim-431 ilarity scores Zhu et al. (2021). However, the frequent encoding and decoding processes in LLMs



Figure 2: Special example for illustrating the effectiveness of invariance loss.

severely hinder efficiency Zhu et al. (2021). To address the trade-off between effectiveness and efficiency, we propose invar-retrieval in our architecture, enabling the model to efficiently retrieve the most relevant documents without introducing variance.

460 Retrieval-augmented Language Model: Currently, retrieval-augmented language models have 461 proven effective in answering questions by leveraging external information through the integration 462 of novel retrievers and LLMs Zhu et al. (2021). However, the architectural gap between retrieval and 463 generation continues to hinder unified optimization across the entire retrieval-augmented generation 464 system Wang et al. (2023). To address the isolation between retrieval and generation, a novel archi-465 tecture called RA-DIT was introduced Lin et al. (2023b). By aligning retriever scoring with LSR 466 scoring Shi et al. (2023), it has been shown to deliver state-of-the-art performance across various 467 tasks. However, it still employs dense retrievers like DRAGON+ Lin et al. (2023a) in the retrieval stage, which fails to eliminate the problem at its source and introduces inefficiencies throughout 468 the process. Correspondingly, we introduce a representation learning method and invariance loss in 469 our Invar-RAG architecture, which partially addresses these issues and explores a novel approach to 470 using a single LLM for multiple roles within the RAG system. 471

472 473

455 456 457

458

459

### 5 CONCLUSION

474 475

476 In this paper, we analyze the challenges and problems of current methods to apply the large language 477 model as a retriever in the RAG system and propose a novel framework, Invar-RAG, to address these 478 challenges. We introduce an LLM-aligned retrieval method, incorporating a well-designed represen-479 tation learning approach to align coarse query-document pairs with the LLM's representation space, 480 allowing our architecture to leverage the extensive parametric knowledge of the LLM to compute 481 relevance scores. Additionally, to address retrieval variance, we propose invariance loss, building 482 on our initial KL-divergence loss, during the retrieval stage to reduce the impact of irrelevant doc-483 uments. Finally, we perform additional fine-tuning on the same LLM for the answer-generation task, enabling our architecture to better utilize the retrieved information and provide more accu-484 rate predictions. Extensive experiments on three open-domain question-answering datasets confirm 485 Invar-RAG's superiority and validate the effectiveness of each module.

## 486 REFERENCES

523

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, et al. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268*, 2016.
- Hangbo Bao, Li Dong, Wenhui Wang, Nan Yang, Songhao Piao, and Furu Wei. Fine-tuning pre trained transformer encoders for sequence-to-sequence learning. *International Journal of Machine Learning and Cybernetics*, 15(5):1711–1728, 2024.
- Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. Rq-rag: Learning to refine queries for retrieval augmented generation. *arXiv preprint arXiv:2404.00610*, 2024.
- Xiaoyang Chen, Ben He, Hongyu Lin, Xianpei Han, Tianshu Wang, Boxi Cao, Le Sun, and Yingfei
  Sun. Spiral of silences: How is large language model killing information retrieval?–a case study on open domain question answering. *arXiv preprint arXiv:2404.10496*, 2024.
- Daixuan Cheng, Shaohan Huang, and Furu Wei. Adapting large language models via reading com prehension. In *The Twelfth International Conference on Learning Representations*, 2023.
- Susan Dumais, Michele Banko, Eric Brill, Jimmy Lin, and Andrew Ng. Web question answering: Is more always better? In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 291–298, 2002.
- Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 6491–6501, 2024.
- Yunlong He, Jiliang Tang, Hua Ouyang, Changsung Kang, Dawei Yin, and Yi Chang. Learning to
   rewrite queries. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, pp. 1443–1452, 2016.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand
   Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning.
   *arXiv preprint arXiv:2112.09118*, 2021.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi
   Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*, 2020.
- Taehyeon Kim, Jaehoon Oh, NakYil Kim, Sangwook Cho, and Se-Young Yun. Comparing kullback-leibler divergence and mean squared error loss in knowledge distillation. *arXiv preprint arXiv:2105.08919*, 2021.
- 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 generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474, 2020.
- Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yuyao Zhang, Peitian Zhang, Yutao Zhu, and Zhicheng Dou.
  From matching to generation: A survey on generative information retrieval. *arXiv preprint arXiv:2404.14851*, 2024.

543

565

567

569

570

540	Sheng-Chieh Lin, Akari Asai, Minghan Li, Barlas Oguz, Jimmy Lin, Yashar Mehdad, Wen-tau Yih,
541	and Xilun Chen. How to train your dragon: Diverse augmentation towards generalizable dense
542	retrieval. arXiv preprint arXiv:2302.07452, 2023a.

- Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Rich James, Pedro Rodriguez, 544 Jacob Kahn, Gergely Szilvasy, Mike Lewis, et al. Ra-dit: Retrieval-augmented dual instruction tuning. arXiv preprint arXiv:2310.01352, 2023b. 546
- 547 Shuming Ma, Dongdong Zhang, and Ming Zhou. A simple and effective unified encoder for 548 document-level machine translation. In Proceedings of the 58th annual meeting of the association for computational linguistics, pp. 3505–3511, 2020. 549
- 550 Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. Fine-tuning llama for multi-stage 551 text retrieval. In Proceedings of the 47th International ACM SIGIR Conference on Research and 552 Development in Information Retrieval, pp. 2421–2425, 2024. 553
- Yuren Mao, Xuemei Dong, Wenyi Xu, Yunjun Gao, Bin Wei, and Ying Zhang. Fit-rag: Black-box 554 rag with factual information and token reduction. arXiv preprint arXiv:2403.14374, 2024. 555
- 556 Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. Large language models: A survey. arXiv preprint arXiv:2402.06196, 558 2024. 559
- Thong Nguyen and Andrew Yates. Generative retrieval as dense retrieval. arXiv preprint 560 arXiv:2306.11397, 2023. 561
- 562 Weizhen Qi, Yu Yan, Yeyun Gong, Dayiheng Liu, Nan Duan, Jiusheng Chen, Ruofei Zhang, and 563 Ming Zhou. Prophetnet: Predicting future n-gram for sequence-to-sequence pre-training. arXiv 564 preprint arXiv:2001.04063, 2020.
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and 566 Yoav Shoham. In-context retrieval-augmented language models. Transactions of the Association for Computational Linguistics, 11:1316–1331, 2023. 568
  - Vipula Rawte, Amit Sheth, and Amitava Das. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922*, 2023.
- 571 Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettle-572 moyer, and Wen-tau Yih. Replug: Retrieval-augmented black-box language models. arXiv 573 preprint arXiv:2301.12652, 2023. 574
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 575 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and 576 efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023. 577
- 578 Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi 579 Yao, Wenyang Gao, Xuming Hu, Zehan Qi, et al. Survey on factuality in large language models: 580 Knowledge, retrieval and domain-specificity. arXiv preprint arXiv:2310.07521, 2023.
- 581 Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-582 attention distillation for task-agnostic compression of pre-trained transformers. Advances in Neu-583 ral Information Processing Systems, 33:5776–5788, 2020. 584
- 585 Krist Wongsuphasawat, Catherine Plaisant, Meirav Taieb-Maimon, and Ben Shneiderman. Querying 586 event sequences by exact match or similarity search: Design and empirical evaluation. *Interacting* with computers, 24(2):55–68, 2012.
- 588 Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. C-pack: 589 Packed resources for general chinese embeddings. In Proceedings of the 47th International ACM 590 SIGIR Conference on Research and Development in Information Retrieval, pp. 641–649, 2024. 591
- Yi Yang, Wen-tau Yih, and Christopher Meek. Wikiqa: A challenge dataset for open-domain ques-592 tion answering. In Proceedings of the 2015 conference on empirical methods in natural language processing, pp. 2013-2018, 2015.

op. 82–86,
nything to
oica, and v preprint
pretrained 024.
eng Chua. 1g. <i>arXiv</i>