# Bridging Relevance and Reasoning: Rationale Distillation in Retrieval-Augmented Generation

Anonymous ACL submission

#### Abstract

The reranker and generator are two critical 002 components in the Retrieval-Augmented Generation (i.e., RAG) pipeline, responsible for ranking relevant documents and generating responses. However, due to differences in pretraining data and objectives, there is an inevitable misalignment between the documents 007 ranked as relevant by the reranker and those required by the generator to support queryspecific answers. To bridge this gap, we 011 propose RADIO, a novel and practical preference alignment framework with RAtionale DIstillatiOn. Specifically, we first propose a ra-013 tionale extraction method that leverages the reasoning capabilities of Large Language Models (LLMs) to extract the rationales necessary for answering a query. Subsequently, a rationale-017 based alignment process is designed to rerank 019 documents based on the extracted rationales and fine-tune the reranker to better align the preferences. Extensive experiments conducted on two tasks across three datasets demonstrate the effectiveness and transferability of our approach. Our code is released online<sup>1</sup>.

## 1 Introduction

027

035

Large Language Models (LLMs), pretrained on massive datasets, have demonstrated exceptional reasoning and text generation capabilities, as evidenced by prior research (Zhao et al., 2023). These models also adhere to the scaling laws, exhibiting improvements in performance and intelligence as the number of parameters increases (Kaplan et al., 2020). Retrieval-augmented generation (RAG) builds upon these capabilities by integrating information retrieval mechanisms with generative models, such as LLMs. This approach not only mitigates the problem of hallucination in text generation but also enhances the system's adaptability to dynamically evolving information needs, making it a robust solution for tasks requiring both accuracy and contextual relevance (Gao et al., 2023). 040

041

042

045

046

047

048

051

052

054

056

057

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

079

However, RAG pipelines typically assemble components (e.g., the reranker and generator (Fan et al., 2024)) that have been pretrained separately. Due to differences in their pretraining data and optimization objectives, these components often exhibit varying preferences, which can impact the overall effectiveness of the system. Specifically, pretrained rerankers (Xiao et al., 2023) excel at evaluating the relevance between queries and documents. However, the documents identified as "relevant" under this criterion may not provide the necessary support for reasoning to derive an accurate answer to the query. Bridging this gap between the reranker's relevance measurement and the generator's reasoning requirements presents a significant challenge that must be addressed to improve the RAG pipeline's performance.

Recent studies try to address this gap by training a bridge model (Ke et al., 2024), using LLMbased scores (Zhang et al., 2024) or combining both LLM-based and retrieval-based scores (Dong et al., 2024) to fine-tune RAG components. Additionally, while some other methods are not explicitly designed for this problem, they can indirectly contribute to bridging the gap. These approaches can use response quality (Ma et al., 2023) or perplexity distillation (Izacard et al., 2023; Shi et al., 2023) as signals to fine-tune the reranker. Despite showing promise, these approaches face critical limitations: their alignment signals rely solely on the surfacelevel connection between the query/answer and document, failing to capture the deeper reasoning processes or more complex relationships involved.

To address the above limitation, we propose RA-DIO, a novel and practical preference alignment framework with rationale distillation in RAG. RA-DIO leverages rationale as a signal to bridge the reranker's relevance measurement with the generator's reasoning requirements for response gen-

<sup>&</sup>lt;sup>1</sup>https://anonymous.4open.science/r/RADIO-9F25

eration. First, to efficiently extract the rationales needed to answer a query, we use the query and its ground truth answer as context and generate the rationales with LLMs. Second, to mitigate the preference misalignment between the reranker and generator while ensuring the solution remains practical, we rerank the documents based on the extracted rationales and fine-tune the reranker. This step distillates rationales from generators to rerankers, and aligns the reranker with the generator's information needs for answering the query effectively.

081

087

094

100

101

104

105

106

107

108

110

111

112

113

114

115

116

117

RADIO effectively addresses the preference inconsistency between RAG components by first generating a comprehensive rationale and then fine-tuning the reranker based on the extracted rationale. This approach considers the deeper reasoning behind answers. We evaluate RADIO on two tasks across three datasets: Open-domain QA (NQ (Kwiatkowski et al., 2019) and TriviaQA (Joshi et al., 2017)) and Multi-choice questions (MMLU (Hendrycks et al., 2020)). The results validate the superiority of our method compared to other state-of-the-art baselines. Our contributions can be summarized as follows:

- We propose RADIO, a novel and practical framework designed to address the preference misalignment of different components in RAG pipelines.
  - We introduce rationale distillation within the RAG framework, which is an effective approach that leverages explicit textual rationales as signals to align the preferences of different components in RAG.
  - Extensive experiments are conducted on two tasks across three datasets to demonstrate the effectiveness and transferability of RADIO.

#### 2 Related Work

#### 2.1 Retrieval-Augmented Generation

Large language models (LLMs) have demon-118 strated groundbreaking performance across numer-119 ous NLP tasks but still face challenges such as hallucination and outdated knowledge (Gao et al., 121 2023). To address these issues, retrieval-augmented 122 generation (RAG) has been introduced (Fan et al., 123 2024). RAG retrieves relevant information from ex-124 125 ternal knowledge bases and incorporates it as contextual input to the generator (LLM), enhancing the 126 accuracy and reliability of the generated responses. 127 The typical RAG pipeline can be divided into several key components: query rewriter, retriever, 129

reranker, and generator. The query rewriter (Wang et al., 2023) modifies and expands the original query to improve retrieval recall, ensuring more relevant documents are retrieved. The retriever (Chen et al., 2024) fetches relevant documents based on the query. Dense retrievers generally outperform sparse retrievers in this step. To integrate contextual information more effectively and identify documents more relevant to the query, rerankers (Moreira et al., 2024a) with larger models and greater complexity are introduced to reorder the retrieved documents compared to retrievers. Finally, the generator-usually a powerful LLM such as GPT-4 (Achiam et al., 2023) or Llama (Touvron et al., 2023)—uses the query and the top-k documents from the reranker to generate the final response. In this work, we address the issue of preference misalignment among different components within the RAG pipeline. We aim to leverage rationale as a signal to align these preferences and enhance the overall performance of the RAG system.

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

#### 2.2 Preference Alignment

To further improve LLMs, preference alignment 152 is often performed after the initial pretraining 153 phase (Jiang et al., 2024). Approaches such 154 RLHF (Ouyang et al., 2022) and DPO (Rafailov 155 et al., 2024) are proposed to align the output 156 of LLMs more closely with human preferences. 157 DPO transforms tasks into classification problems, 158 achieving high computational efficiency and strong 159 performance. In the context of RAG, several 160 works (Ke et al., 2024) can be transformed to ad-161 dress the challenge of preference alignment be-162 tween RAG components. REPLUG (Shi et al., 163 2023) improves RAG pipelines involving black-164 box LLMs by using the probability of the LLM 165 generating the correct answer as a signal to deter-166 mine document importance. Similarly, RRR (Ma 167 et al., 2023) uses metrics based on the quality 168 of the LLM's generated response as a signal to 169 evaluate a document's utility. On the other hand, 170 ARL2 (Zhang et al., 2024) prompts LLMs to gen-171 erate self-guided relevance labels for fine-tuning 172 retriever, and DPA-RAG (Dong et al., 2024) intro-173 duces a bidirectional alignment strategy to mitigate 174 preference inconsistencies in RAG pipelines. In 175 this work, we focus on optimizing the reranker 176 within RAG. Our goal is to enable the reranker 177 to effectively identify supportive documents well-178 suited for the generator, whether black-box or open-179 source, facilitating the production of accurate out-180

275

276

181puts. In addition, BGM (Ke et al., 2024) trains182a bridge model between retriever and LLMs to183transform the retrieved information into the format184LLM's prefer. Our method is theoretically compat-185ible with BGM and could be combined to jointly186enhance performance without any conflict.

## 3 Methodology

187

188

190

191

192

193

194

195

196

197

198

199

200

206

207

208

209

210

212

213

214

215

216

217

218

219

221

226

In this section, we detail rationale distillation in RAG. Specifically, we first demonstrate the task definition of RAG in Section 3.1. Then we give an overview of our proposed framework in Section 3.2, introduce the rationale extraction method in Section 3.3, and detail the rationale-based alignment in Section 3.4 and optimization in Section 3.5.

#### 3.1 Task Definition

To address the hallucination problem and enhance adaptability to dynamic information of LLMs, RAG systems have been proposed. These systems enhance generative models by introducing additional contextual information retrieved based on a given query q. Specifically, when a query q is input into the RAG pipeline, the retriever  $\mathcal{R}_{\text{retriever}}$  first retrieves relevant documents by calculating similarity scores and top- $k_1$  selection. The process can be formalized as follows:

$$\mathcal{D}_{\text{retriever}} = \{ \boldsymbol{d}_i \mid \boldsymbol{d}_i \in \text{Top-}k_1(\text{score}_{\text{retriever}}(q, d)) \}$$
(1)

where q and d are the query and document, score<sub>retriever</sub> denotes the score function in retriever,  $d_i$  means the *i*-th document in corpus and  $\mathcal{D}_{retriever}$ is the documents set output by retriever  $\mathcal{R}_{retriever}$ .

To eliminate contextually irrelevant noise and provide more precise contextual information, the initially filtered documents will be further reranked by the reranker:

$$\mathcal{D}_{\text{reranker}} = \{ \boldsymbol{d}_j \mid \boldsymbol{d}_j \in \text{Top-}k_2(\text{score}_{\text{reranker}}(q, d)) \}$$
(2)

where  $\mathcal{D}_{\text{reranker}}$  is the documents selected by reranker, score<sub>reranker</sub> denotes the score function in reranker. Note that the *j*th document  $d_j$  is in  $\mathcal{D}_{\text{retriever}}$  (i.e.,  $d_j \in \mathcal{D}_{\text{retriever}}$ ) and  $k_2$  is the number of documents selected by reranker, which is smaller than  $k_1$  used by the retriever.

Finally, the documents filtered by the reranker, along with the original query, will be fed into the generator as contextual information to help generate the final response:

$$\hat{y} = \mathcal{G}(q, \mathcal{D}_{\text{reranker}})$$
 (3)

where  $\hat{y}$  is the generated response and  $\mathcal{G}$  denotes the generator.

It is worth noting that the documents selected by the reranker directly influence the generator's input. Therefore, in this work, we aim to align the preferences of the reranker and generator to enhance their consistency. This alignment improves the overall accuracy of the RAG system's responses.

#### 3.2 Framework Overview

The overview of RADIO is depicted in Figure 1. RADIO is consisted of two phases: rationale extraction (Figure 1(a)) and rationale-based alignment (Figure 1(b)).

In the rationale extraction process, we combine the query with its ground truth answer and input them into LLMs to generate precise rationales. Using the correct answer as context in the prompt improves the accuracy of the LLM's rationale generation, ensuring that the generated rationale closely aligns with the requirements for deriving the correct answer.

In the rationale-based alignment process, our goal is to use the generated rationale to guide the reranker, enabling it to select documents that better support the generator in answering the query. Specifically, we leverage the generated rationale as a signal to rerank documents. The reranked documents will be used to fine-tune the reranker, addressing the preference misalignment between the reranker and the generator. By aligning these components, the process ensures that the selected documents are not only contextually relevant but also optimally supportive for the generator's reasoning and response generation.

#### 3.3 Rationale Extraction

Rationales are critical components of LLM reasoning processes and have been shown to significantly enhance the accuracy of LLM-generated responses. This perspective is supported by existing works such as Chain-of-Thoughts (CoT (Wei et al., 2022)) and O1 (Zhong et al., 2024). Existing work (Shi et al., 2023; Ma et al., 2023; Dong et al., 2024) has primarily focused on the initial relationships between queries and documents or indirect relationships between answers and documents, while overlooking rationales, a crucial intermediary component in the reasoning process. Motivated by this, we aim to extract rationales as signals to align the preferences of different components in the RAG pipeline.



Figure 1: Overview of RADIO.

To accurately extract the rationale necessary for answering a query and deriving the correct answer, we combine the query with the ground truth answer as contextual information of LLMs, as shown in Figure 1(a). The prompt template used for this process is as follows: "You are a professional QA assistant. Given a question and the ground truth answer, you can output the rationale why the ground truth answer is correct. Question: {question}. Answer: {answer}. Rationale: ". The generation process can be formalized as:

277

279

282

283

287

290

291

294

297

301

305

307

313

$$r = \text{LLM}(q, a) \tag{4}$$

where q and a are the query and answer, r deontes the generated rationale.

By doing so, we effectively bridge the gap between the query and the answer by generating the necessary rationale. This rationale accurately supports the reasoning process required to derive the correct answer from the query.

#### 3.4 Rationale-based Alignment

Given the extracted rationale r, a key challenge lies in effectively and efficiently utilizing it to improve preference consistency within the RAG pipeline. In this section, we propose a rationale-based alignment approach, where the rationale serves as a signal to fine-tune the reranker. This enables the reranker to identify and prioritize supportive documents that facilitate the generator in producing accurate responses. Specifically, we first use the retriever  $\mathcal{R}_{retriever}$  to retrieve  $k_1$  relevant documents based on the query:

$$\mathcal{D}_{\text{retriever}} = \mathcal{R}_{\text{retriever}}(q, \mathcal{C}) \tag{5}$$

where  $\mathcal{D}_{retriever}$  is the document set retrieved by  $\mathcal{R}_{retriever}$  based on query q and corpus C, and  $|\mathcal{D}_{retriever}| = k_1$ .

Next, to facilitate the comparison of similarity between different documents and the rationale, we

use a text encoder to convert both the documents and the rationale into dense vectors.

e

$$p_i^{\text{document}} = \text{Encoder}(\boldsymbol{d}_i)$$
 (6)

$$e^{\text{rationale}} = \text{Encoder}(r)$$
 (7) 3

314

315

316

318

319

320

321

322

323

324

325

327

328

329

330

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

where  $e_i^{\text{document}}$  and  $e^{\text{rationale}}$  denote the representations of *i*th document and rationale.  $d_i$  is the *i*th document and *r* represents the extracted rationale.

Then, we calculate the semantic similarity between each document and the rationale. Here, we use cosine similarity, denoted as  $sim(\cdot)$ . The calculated scores indicate the degree to which each document supports generating the correct answer, with higher scores reflecting stronger support. We also linearly interpolate the score of documents with their retrieval score in the retrieval stage by weighted score sum.

$$s_i^{\text{rationale}} = \sin(e_i^{\text{document}}, e^{\text{rationale}})$$
 (8)

$$s_i^{\text{retriever}} = \text{score}_{\text{retriever}}(q, d_i)$$
 (9)

$$\mathbf{s}'_i = \alpha \mathbf{s}^{\text{rationale}}_i + (1 - \alpha) \mathbf{s}^{\text{retriever}}_i$$
 (10)

where  $s_i^{\text{rationale}}$ ,  $s_i^{\text{retriever}}$ ,  $s_i'$  represent the rationale similarity score, retrieval score, and final score for *i*th document. score<sub>retriever</sub>(·) is the score function in retriever and  $\alpha$  is a hyperparameter used for integration. Note that we apply min-max normalization in this work to both rationale score and retrieval score before integration.

Next, we rerank the documents based on their scores. Following the previous sampling method *Top-k shifted by N* (Moreira et al., 2024b), we select the top-ranked document as the positive sample and then shift by n documents and sample m negative samples from the subsequent documents to construct positive-negative pairs for fine-tuning the reranker. This process can be represented as:

$$d_{\text{pos}} = d_i$$
, where  $i = \arg \max_i s'_i$  (11) 34

$$\{\boldsymbol{d}_{\text{neg}}\} = \text{Sample}_m(\{\boldsymbol{d}_i | \text{rank}(\boldsymbol{d}_i) > n\}) \quad (12) \quad 34$$

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

394

395

396

397

where  $d_{pos}$  and  $d_{neg}$  are the sampling positive and negative documents, Sample<sub>m</sub>(·) denotes a sampling operation that selects m negative documents from the set of documents ranked lower than n.

#### 3.5 Optimization

356

358

362

363

370

372

374

375

376

377

Following BGE embedding (Xiao et al., 2023) and QA Ranking Benchmark (Moreira et al., 2024a), we use InfoNCE as our optimization objectives to fine-tune reranker:

$$f(q,d) = \exp(\phi(q,d)/\tau)$$
(13)

$$L = -\log \frac{f(q, d^{+})}{f(q, d^{+}) + \sum_{i=1}^{N} f(q, d_{i}^{-})}$$
(14)

where  $d^+$  and  $d^-$  represent the positive and negative document,  $\tau$  is the temperature parameter, and N denotes the number of negative documents.

## 4 Experiments

## 4.1 Datasets and Metrics

To evaluate RADIO with other methods, we conduct experiments on two tasks across three datasets: **Open-domain QA** (NQ (Kwiatkowski et al., 2019) and TriviaQA (Joshi et al., 2017)) and **Multi-choice questions** (MMLU (Hendrycks et al., 2020)). Deatiled dataset descriptions are given in Appendix A.1. Following previous work (Ma et al., 2023; Shi et al., 2023), we report EM and F1 scores for Open-domain QA datasets and EM for MMLU.

#### 4.2 Baselines

To verify the effectiveness of RADIO, we conduct experiments with the following baseline methods: **Base** (Xiao et al., 2023), **Atlas** (Izacard et al., 2023), **REPLUG** (Shi et al., 2023), **Trainable rewrite-retrieve-read** (**RRR**) (Ma et al., 2023), **ARL2** (Zhang et al., 2024), and **DPA-RAG** (Dong et al., 2024). The detailed introduction of baselines is given in Appendix A.2.

#### 4.3 Backbone Rerankers

385To validate the generality and adaptability of RA-386DIO, we select three different rerankers as the back-387bone models for our experiments: (1) gte-base (Li388et al., 2023): a reranker model proposed by Alibaba389DAMO Academy, with 109M parameters and 768390embedding dimensions. (2) gte-large (Li et al.,3912023): the larger version of gte-base, with 335M392parameters and 1024 embedding dimensions. (3)393bge-reranker-base (Xiao et al., 2023): a powerful

cross-encoder architecture reranker proposed by Beijing Academy of Artificial Intelligence, with 278M parameters.

#### 4.4 Implementation Details

We implement RADIO on FlashRAG (Jin et al., 2024), a Python library for efficient RAG research. In the RAG pipeline, we take e5-base-v2 (Wang et al., 2022) as the retriever, and Meta-Llama-3.1-8B-Instruct (Touvron et al., 2023) as the generator. We sample 20,000 instances from NO and TriviaQA to construct the fine-tuning dataset and finetune rerankers separately. For document sampling, we set the shift n in Top-k shifted by N method as 3, and sample 6 negative samples from the subsequent documents. To ensure a fair comparison, the sampling index is fixed and remains unchanged across methods. In the RAG pipeline, we set the number of documents selected by retriever and reranker (i.e.,  $k_1$  and  $k_2$ ) as 20 and 5. For fine-tuning the reranker, we tune the training epochs from 1 to 5 and the integration hyperparameter  $\alpha$  from 0.0 to 1.0. We use Adam (Kingma, 2014) optimizer with a learning rate 6e-5 and a weight decay of 0.01. The prompts we used in experiments are given in Appendix A.5.

#### 4.5 Main Results

## 4.5.1 Open-domain QA

To evaluate the effectiveness of RADIO and its transferability across different rerankers, we conduct experiments on the NQ and TriviaQA datasets. The results are presented in Table 1. From these results, we can draw the following conclusions:

- Compared to the Base method, most experimental settings achieve better results, demonstrating the necessity of preference alignment within the RAG pipeline.
- Compared to other baseline methods, RA-DIO consistently achieves superior performance across all datasets and reranker backbone configurations. The results validate the effectiveness of using rationales as signals for preference alignment in RAG pipeline.
- On TriviaQA, methods such as RRR and RE-PLUG show performance declines relative to the base method when using rerankers *gte-large* and *bge-reranker-base*. This indicates that these methods are sensitive, limiting their applicability.

Table 1: Overall experiments. "\*" indicates the statistically significant improvements (i.e., two-sided t-test with p < 0.05) over the best baseline. For all metrics, higher is better.  $\Delta$  represents the relative improvement of RADIO over Base method.

	NQ							TriviaQA						
Method	gte-base		gte-large		bge-reranker-base		gte-base		gte-large		bge-reranker-base			
	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1		
Base	0.2931	0.4046	0.2798	0.3935	0.3371	0.4603	0.5449	0.6374	0.5495	0.6414	0.6114	0.7120		
Atlas	0.3338	0.4587	0.3418	0.4677	0.3521	0.4832	0.5823	0.6752	0.6004	0.6972	0.6083	0.7063		
REPLUG	0.3257	0.4484	0.2607	0.3670	0.3427	0.4753	0.5679	0.6578	0.5248	0.6171	0.6032	0.7004		
RRR	0.3374	0.4608	0.3299	0.4578	0.3438	0.4754	0.5801	0.6716	0.5358	0.6237	0.6099	0.7091		
ARL2	0.3413	0.4688	0.3515	0.4804	0.3568	0.4885	0.6079	0.7086	0.6107	0.7120	0.6137	0.7149		
DPA-RAG	0.3391	0.4674	0.3385	0.4710	0.3462	0.4793	0.6080	0.7076	0.6097	0.7119	0.6149	0.7169*		
RADIO (Ours)	0.3512*	0.4790*	0.3565*	0.4850*	0.3665*	0.4917*	0.6084	0.7095*	0.6128*	0.7137*	0.6154	0.7151		
Δ	$\uparrow 19.82\%$	$\uparrow 18.39\%$	↑27.41%	↑23.25%	<b>↑8.72%</b>	↑ <b>6.82%</b>	↑11.55%	↑11.31%	↑11.52%	<u></u> ↑11.27%	↑0.65%	↑0.40%		

Table 2: Experimental results on MMLU. EM is reported as the evaluation metric. The source dataset used to fine-tune rerankers is the Open-domain QA dataset NQ.  $\Delta$  represents the relative improvement of RADIO over Base method.

Method	Humanitie	es Social	STEM	Other	ALL
Base	0.4089	0.6867	0.5147	0.6650	0.5502
Atlas	0.3985	0.6935	0.5074	0.6563	0.5447
REPLUG	0.4102	0.6854	0.5065	0.6590	0.5473
RRR	0.4079	0.6913	0.5116	0.6572	0.5484
ARL2	0.4147	0.7016	0.5106	0.6630	0.5540
DPA-RAG	0.4157	0.701	0.5078	0.6652	0.5541
RADIO (Ours)	0.4172	0.7013	0.5080	0.6717	0.5562
$\Delta$	↑2.03%	$\uparrow 2.13\%$	↓1.30%	$\uparrow 1.01\%$	$\uparrow 1.09\%$

In contrast, RADIO demonstrates robust adaptability to different rerankers, achieving significant performance improvements across all three rerankers.

• As the reranker becomes larger or more powerful (e.g., progressing from *gte-base* to *gte-large* and further to *bge-reranker-base*), the performance ranking of models fine-tuned with RADIO aligns with the reranker's inherent capabilities. This suggests that RADIO's performance gains are sustainable and scalable with stronger rerankers, providing an avenue to further explore the upper performance limits of RAG pipelines.

## 4.5.2 MMLU

441

442

443

444

445

446

447

448

449

450

451

452

453

454

We also conduct experiments on MMLU. Since 455 MMLU is a multiple-choice dataset, we report the 456 EM metric (Ma et al., 2023). Additionally, follow-457 ing previous work (Yu et al., 2023), we fine-tune the 458 reranker using open-domain QA as the source task 459 460 and evaluate its performance on the MMLU dataset. Table 2 shows the results of fine-tuning reranker 461 with NQ dataset. The results of fine-tuning reranker 462 with TriviaQA are given in Appendix 5. We can 463 draw the following conclusions: 464

From the metrics corresponding to the ALL category, RADIO demonstrates consistent improvements over Base. This highlights the effectiveness and transferability of RADIO, as it successfully adapts to multi-choice question tasks even when fine-tuned on the Open-domain QA tasks.

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

- Analyzing the results by question category, RA-DIO shows more significant improvements over the Base method in the Humanities and Social Sciences categories, with average gains of 2.03% and 2.13%, respectively. However, it exhibits a slight negative effect in the STEM category. This may be due to the fine-tuning datasets (NQ and TriviaQA), which are Open-domain QA datasets with distributions more similar to humanities and social sciences but markedly different from STEM subjects.
- Compared to other baseline methods, RADIO achieves top performance in the vast majority of metrics, demonstrating its superiority and state-of-the-art capability.

## 4.6 Transferability Analysis across Generators

We conduct experiments on two datasets of different tasks, NQ and MMLU, using three different generators (*Llama3.1-8b-instruct* (Touvron et al., 2023), *qwen2.5-14b-instruct* (Yang et al., 2024), and *gpt4o-mini* (Achiam et al., 2023)) to validate the transferability of our method, as shown in Table 3. We can find: (1) RADIO maintains its effectiveness across different generators, consistently enhancing the performance of the original RAG pipeline. This demonstrates RADIO's strong transferability with various generators. (2) Comparing different generators reveals that RA-DIO's performance gains are more pronounced

Generator	Method	NQ		MMLU					
Contractor		EM	F1	Humanities	Social	STEM	Other	ALL	
I lama 1_8h_instruct	Base	0.3371	0.4603	0.4089	0.6867	0.5147	0.6650	0.5502	
	RADIO (ours)	0.3665	0.4917	0.4172	0.7013	0.5080	0.6717	0.5562	
awan 2.5.14h instruct	Base	0.3310	0.4484	0.5439	0.8200	0.7206	0.7541	0.6906	
qwell2.3-140-llistiuct	RADIO (ours)	0.3518	0.4753	0.5598	0.8229	0.7250	<b>0.763</b> 1	0.6995	
ant la mini	Base	0.3607	0.4880	0.6485	0.8362	0.6784	0.8005	0.7300	
gpt40-mm	RADIO (ours)	0.3742	0.5086	0.6548	0.8372	0.6768	0.8010	0.7321	

Table 3: Transferability analysis across generators. EM is reported as the metric for MMLU dataset.

Table 4: Ablation study.

Dataset	Metrics	w/o ALL	w/o Retrieval	RADIO
NQ	EM F1	0.3371 0.4603	$\frac{0.3587}{0.4858}$	0.3665 0.4917
MMLU	Humanities (EM) Social (EM) STEM (EM) Other (EM) ALL (EM)	0.4089 0.6867 <b>0.5147</b> 0.665 0.5502	$\begin{array}{r} 0.4168 \\ \hline 0.6981 \\ \hline 0.5109 \\ \hline 0.6666 \\ \hline 0.5548 \end{array}$	0.4172 0.7013 0.508 0.6717 0.5562

with smaller, less capable generators. Specifically, when the generators are Llama3.1, Qwen2.5, and GPT4o-mini, RADIO achieves EM improvements of 8.72%, 6.28%, and 3.74%, respectively, and F1 improvements of 6.82%, 6.00%, and 4.22%. This is because as the generator's capability increases and approaches the upper performance limits of the RAG pipeline, further enhancing the pipeline becomes increasingly challenging, resulting in a smaller improvement.

#### 4.7 Ablation Study

501

502

504

508

510

511

512

513

514

515

516

517

518

521

523

525

529

To explore the specific impact of rationale and retrieval score, we design the following variants: (1) w/o ALL: Base reranker without fine-tuning. Do not introduce rationale or retrieval score. (2) w/o Retrieval: Ranking documents and fine-tuning reranker only based on the rationale scores. (3) RADIO: Fine-tuning reranker based on both rationale and retrieval scores.

Table 4 shows the results of ablation study on NQ and MMLU, where we can derive the following findings: (1) Both the rationale score and retrieval score contribute positively to RADIO's performance, with the rationale score demonstrating a stronger positive impact compared to the retrieval score. (2) RADIO outperforms both w/o ALL and w/o Retrieval, while w/o Retrieval surpasses w/o ALL. This indicates that the rationale score and retrieval score are not conflicting but rather com-



Figure 2: Hyperparameter analysis on NQ and MMLU.

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

plementary. Their integration provides a more robust signal for document ranking, which effectively aids in fine-tuning the reranker. (3) In the STEM category of the MMLU dataset, **w/o ALL** outperforms both **w/o Retrieval** and RADIO. This could be attributed to the fact that the training dataset (NQ) contains questions with distributions more similar to humanities and social sciences, leading to trends in the Humanities, Social Sciences, and Other categories that differ from STEM category.

#### 4.8 Hyperparameter Analysis

Figure 2 visualizes the performance of RADIO across different integration coefficients  $\alpha$  on NQ and MMLU. The x-axis represents the integration coefficients  $\alpha$  and the y-axis represents the evaluation metrics EM (red) or F1 (blue). For the MMLU dataset, we present results for two representative categories: Humanities and STEM. The trends for other categories align with those observed in Humanities. Complete experimental results are provided in Appendix A.4 for reference. From the figure, we can draw the following conclusions: (1) As  $\alpha$  increases, the metrics on the NQ dataset and most

categories of the MMLU dataset exhibit a trend of 553 first rising and then falling, with the optimal range 554 for  $\alpha$  being between 0.3 and 0.7. This demonstrates 555 the complementary nature of rationale scores and retrieval scores, which together form an optimal signal for preference alignment. (2) When  $\alpha = 0$ , the RAG performance is suboptimal because docu-559 ment scoring relies entirely on retrieval scores, focusing solely on query-document relevance while ignoring whether the document supports the gen-562 erator in answering the query. Conversely, when  $\alpha = 1$ , the performance is still not optimal, as it 564 completely disregards retrieval relevance, leading 565 to a mismatch between the fine-tuning dataset and the training data, which negatively affects model 567 performance. (3) The trends for STEM differ from those of other MMLU categories, showing an opposite pattern. This is likely due to the significant distributional differences between STEM and Hu-571 manities/Social Sciences, resulting in a "seesaw effect" as observed in the figure. This phenomenon is 573 also reflected in REPLUG (Shi et al., 2023), where the improvement in the STEM category is weaker compared to other categories. 576

## 4.9 Case Study

To intuitively illustrate the effectiveness of RADIO, we select examples from the NQ dataset to compare the documents reranked by RADIO with those reranked by a novel baseline, DPA-RAG, as well as 581 their responses. In Figure 3 Example 1, the query asks, "Which state is the richest state in Nigeria?" RADIO successfully ranks information about La-585 gos State's economic and financial status, which relates to the correct answer, among the top-3 doc-586 uments. In contrast, DPA-RAG fails to identify 587 documents relevant to answering the query, and cannot provide a valid response. In Figure 3 Exam-589 ple 2, the query is, "Who is the highest-selling R&B artist of all time?" RADIO prioritizes documents 591 containing information about the correct answer, 592 Michael Jackson, and effectively highlights key terms such as "R&B" and "best-selling." However, 594 DPA-RAG misinterprets the query's constraints, retrieving documents that either overlook the R&B artist specification or fail to consider the time span, 598 resulting in an incorrect response. These examples demonstrate that RADIO enhances RAG by providing a more efficient and accurate reranking. It selects contextually appropriate documents, enabling the generator to infer correct answers. 602



Figure 3: Case study on NQ dataset.

603

604

605

606

607

608

609

610

611

612

613

614

615

# 5 Conclusion

In this paper, we propose a novel and practical preference alignment framework, RADIO, with rationale distillation in retrieval-augmented generation. First, we introduce a rationale extraction method to extract the rationales necessary for answering queries with LLMs. Next, a rationale-based alignment is proposed to rerank documents based on extracted rationales and fine-tune rerankers. Extensive experiments on two tasks across three datasets are conducted to validate the effectiveness of our proposed method against state-of-the-art baselines and demonstrate its strong transferability.

#### 6 Limitations

616

633

634

635

647

651

657

658

662

617First, compared to other methods, our approach618RADIO requires additional time in the rationale619extraction stage to generate rationales. Since dif-620ferent samples are independent of one another, we621can reduce generation time by employing paral-622lel processing to mitigate this issue. Secondly,623the MMLU experimental results reveal that the624composition of fine-tuning datasets can affect RA-625DIO's effectiveness. This issue can be addressed626by designing task-specific fine-tuning datasets for627different downstream tasks or large-scale general628fine-tuning datasets.

#### References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. *Preprint*, arXiv:2402.03216.
- Guanting Dong, Yutao Zhu, Chenghao Zhang, Zechen Wang, Zhicheng Dou, and Ji-Rong Wen. 2024. Understand what llm needs: Dual preference alignment for retrieval-augmented generation. *arXiv preprint arXiv:2406.18676*.
- Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. 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, pages 6491– 6501.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:2312.10997.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
  2020. Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot learning with retrieval augmented language models. *Journal of Machine Learning Research*, 24(251):1–43.

Ruili Jiang, Kehai Chen, Xuefeng Bai, Zhixuan He, Juntao Li, Muyun Yang, Tiejun Zhao, Liqiang Nie, and Min Zhang. 2024. A survey on human preference learning for large language models. *arXiv preprint arXiv:2406.11191*. 667

668

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

- Jiajie Jin, Yutao Zhu, Xinyu Yang, Chenghao Zhang, and Zhicheng Dou. 2024. Flashrag: A modular toolkit for efficient retrieval-augmented generation research. *arXiv preprint arXiv:2405.13576*.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Zixuan Ke, Weize Kong, Cheng Li, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. 2024. Bridging the preference gap between retrievers and llms. *arXiv preprint arXiv:2401.06954*.
- Diederik P Kingma. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453– 466.
- Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. 2023. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*.
- Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. Query rewriting for retrievalaugmented large language models. *arXiv preprint arXiv:2305.14283*.
- Gabriel de Souza P Moreira, Ronay Ak, Benedikt Schifferer, Mengyao Xu, Radek Osmulski, and Even Oldridge. 2024a. Enhancing q&a text retrieval with ranking models: Benchmarking, fine-tuning and deploying rerankers for rag. *arXiv preprint arXiv:2409.07691*.
- Gabriel de Souza P Moreira, Radek Osmulski, Mengyao Xu, Ronay Ak, Benedikt Schifferer, and Even Oldridge. 2024b. Nv-retriever: Improving text embedding models with effective hard-negative mining. *arXiv preprint arXiv:2407.15831*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al.

- 722 727 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748

- 754 755 757 759 761

- 772

774

775

776

777

778

767

768 769 770 771

2022. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730–27744.

- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. Advances in Neural Information Processing Systems, 36.
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2023. Replug: Retrievalaugmented black-box language models. arXiv preprint arXiv:2301.12652.
- Devendra Singh, Siva Reddy, Will Hamilton, Chris Dyer, and Dani Yogatama. 2021. End-to-end training of multi-document reader and retriever for opendomain question answering. Advances in Neural Information Processing Systems, 34:25968–25981.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weaklysupervised contrastive pre-training. arXiv preprint arXiv:2212.03533.
- Liang Wang, Nan Yang, and Furu Wei. 2023. Query2doc: Query expansion with large language models. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 9414-9423.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837.
- Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023. C-pack: Packaged resources to advance general chinese embedding. Preprint, arXiv:2309.07597.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu

Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 technical report. arXiv preprint arXiv:2407.10671.

779

780

781

782

783

784

785

787

789

790

791

792

793

794

795

796

797

798

- Zichun Yu, Chenyan Xiong, Shi Yu, and Zhiyuan Liu. 2023. Augmentation-adapted retriever improves generalization of language models as generic plug-in. arXiv preprint arXiv:2305.17331.
- Lingxi Zhang, Yue Yu, Kuan Wang, and Chao Zhang. 2024. Arl2: Aligning retrievers for black-box large language models via self-guided adaptive relevance labeling. arXiv preprint arXiv:2402.13542.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. arXiv preprint arXiv:2303.18223.
- Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, Shizhe Liang, Zihao Wu, Yanjun Lyu, Peng Shu, Xiaowei Yu, et al. 2024. Evaluation of openai o1: Opportunities and challenges of agi. arXiv preprint arXiv:2409.18486.

10

803

810

811

812

814

815

816

817

818

819

821

822

823

825

826

827

832

834

839

840

841

842

843

# A Appendix

# A.1 Dataset Descriptions

The detailed descriptions of baselines are given as follows:

- Natural Questions (NQ): NQ contains real user questions compiled from Google search, with corresponding answers identified from Wikipedia by human annotators.
- **TriviaQA:** TriviaQA dataset comprises trivia questions paired with answer annotations and supporting evidence documents, such as web pages and Wikipedia articles. It is designed to assess a model's ability to retrieve and comprehend textual evidence for open-domain question answering.
  - Massive Multitask Language Understanding (MMLU): MMLU is a comprehensive evaluation dataset comprising 57 categories of questions, which are grouped into four broad domains: Humanities, Social Sciences, STEM, and Other. In this paper, we report evaluation metrics based on these categories.

## A.2 Baselines

Following is the introduction of baselines:

- **Base (Xiao et al., 2023):** The reranker model is used off-the-shelf without any fine-tuning.
- Atlas (Izacard et al., 2023): A pretrained retrieval-augmented language model designed for knowledge intensive task. We choose the EMDR<sup>2</sup> (Singh et al., 2021) as the reward to rerank documents and fine-tune rerankers.
- **REPLUG (Shi et al., 2023):** REPLUG seeks to fine-tune the retriever to enhance RAG pipelines that include black-box LLMs. It achieves this by using the query and document as contextual inputs and leveraging the probability of the LLM generating the correct answer as the importance score. This idea is also reflected in the PDist method in Atlas (Izacard et al., 2023).
- Trainable rewrite-retrieve-read (RRR) (Ma et al., 2023): RRR optimizes the query rewriter using the evaluation metrics of the final RAG output as a reward, which is used to fine-tune rerankers in our pipeline, enhancing the overall effectiveness of RAG.

• ARL2 (Zhang et al., 2024): ARL2 introduces a method to use LLMs as supervisor to generate self-guided relevance labels for fine-tuning retriever.

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

• **DPA-RAG** (**Dong et al., 2024**): DPA-RAG proposes a knowledge preference pipeline to dual-align rerankers and generators in RAG. It combines document importance from the LLM's perspective with the importance determined during the retrieval stage.

# A.3 More Results on MMLU

In this section, we give the complete experimental results on MMLU. Specifically, we use the NQ dataset and TriviaQA dataset as source dataset to fine-tune rerankers and evaluate them in MMLU.

# A.4 More Results on Hyperparameter Analysis

Figure 4 illustrates the trend of RADIO's performance across various MMLU categories as the integration coefficient  $\alpha$  increases. For Humanities, Social Sciences, and Other, the trends are consistent: as  $\alpha$  increases, the performance metrics first improve and then decline. However, the STEM category shows a unique pattern, with metrics initially decreasing as  $\alpha$  grows, followed by an improvement. This divergence may be attributed to the fine-tuning dataset (NQ), which shares a closer distribution with Humanities, Social Sciences, and Other categories, while differing significantly from the STEM category.

# A.5 Prompts

In this section, we detail the prompts we used in the experiments. For Open-domain QA datasets NQ and TriviaQA, the prompts are as shown in Table 6. For Multi-choice dataset MMLU, the prompts are shown in Table 7.



Figure 4: More results of hyperparameter analysis on MMLU.

Table 5: Experimental results on MMLU. EM is reported as the metric. The source datasets used to fine-tune rerankers are Open-domain QA datasets NQ and TriviaQA.

Method		MMLU (	Source Dat	taset NQ)		MMLU (Source Dataset TriviaQA)				
	Humanitie	es Social	STEM	Other	ALL	Humanitie	es Social	STEM	Other	ALL
Base	0.4089	0.6867	0.5147	0.6650	0.5502	0.4089	0.6867	0.5147	0.6650	0.5502
Atlas	0.3985	0.6935	0.5074	0.6563	0.5447	0.3966	0.6822	0.4868	0.654	0.5364
REPLUG	0.4102	0.6854	0.5065	0.6590	0.5473	0.3977	0.6744	0.4821	0.6466	0.5323
RRR	0.4079	0.6913	0.5116	0.6572	0.5484	0.4132	0.6926	0.5005	0.6501	0.5464
ARL2	0.4147	0.7016	0.5106	0.6630	0.5540	0.4012	0.6951	0.5011	0.6639	0.5462
DPA-RAG	0.4157	0.701	0.5078	0.6652	0.5541	0.4189	0.6932	0.5062	0.6746	0.5552
RADIO (Ours)	0.4172	0.7013	0.5080	0.6717	0.5562	0.4230	0.6942	0.5090	0.6678	0.5559

Table 6: Prompts for Open-domain datasets.

System Prompt: Answer the question based on
the given document. Only give me the answer
and do not output any other words. The
following are given documents.
{reference}
User Prompt:
Question: {question}
Answer:

Table 7: Prompts for Multi-choice datasets.

System Prompt: Answer the question based on the given document. Only give me the option (A/B/C/D) and do not output any other words. The following are given documents. {reference} User Prompt: Question: {question} Answer: