

# Rethinking Retrieval-Augmented Generation as a Cooperative Decision-Making Problem

Anonymous ACL submission

## Abstract

Retrieval-Augmented Generation (RAG) has demonstrated strong effectiveness in knowledge-intensive tasks by grounding language generation in external evidence. Despite its success, many existing RAG systems are built based on a ranking-centric, asymmetric dependency paradigm, where the generation quality of the generator is highly dependent on reranking results of the reranker. To overcome this limitation, we reformulate RAG as a cooperative multi-agent decision-making problem and propose Cooperative Retrieval-Augmented Generation (CoRAG), a framework in which the reranker and the generator act as peer decision-makers rather than being connected through an asymmetric dependency pipeline. By jointly optimizing their behaviors toward a shared task objective, the reranker and generator are encouraged to cooperate, ensuring that document reranking and generation work in concert to improve the final response. Experimental results demonstrate good generalization and improved generation stability of CoRAG, even when the model is trained on only around 10K PopQA samples. Our model released in <https://anonymous.4open.science/r/CoRAG-D63F>

## 1 Introduction

In recent years, Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Asai et al., 2024; Gao et al., 2023b; Zhao et al., 2024) has emerged as an important paradigm for enhancing the factuality and knowledge coverage of large language models (LLM) (Gao et al., 2023b). A typical RAG system consists of two core components: a retriever and a generator (Lewis et al., 2020; Oche et al.). Given a query, the retriever retrieves candidate documents from a large external corpus and further rerank them via reranker. The generator then conditions on the query and the reranked documents to produce the final response. By explicitly incorporating

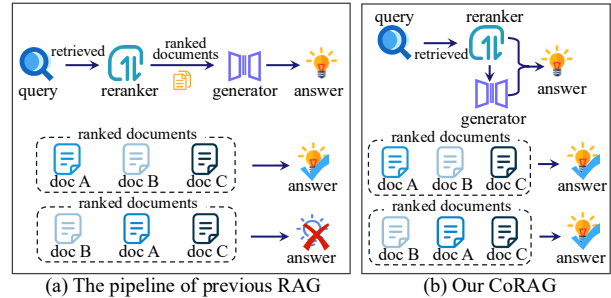


Figure 1: Comparison with previous works. Previous works assume an asymmetric dependency between the reranker and the generator, whereas CoRAG models them as cooperative agents in a multi-agent reinforcement learning framework.

external knowledge during generation, RAG effectively mitigates hallucinations and achieves strong performance on open-domain question answering tasks (Lewis et al., 2020; Asai et al., 2024).

As the reranker plays a critical role in shaping the document context provided to the generator (Lewis et al., 2020; Sharma, 2025), a growing body of recent work has focused on the design and optimization of rerankers and generators (Gao et al., 2023b; Asai et al., 2024; Sun et al., 2025; Shen et al., 2023; Jia et al., 2025). However, most existing RAG methods still adopt a ranking-centric, asymmetric-dependency paradigm (Figure 1(a)), where the reranker produces a fixed document ordering and the generator performs generation conditioned on these top-ranked documents. This design tightly couples generation with reranking decisions, making the generator highly sensitive to the reranking results. As shown in Figure 1(a), if a suboptimal reranker misranks a less relevant document at the top, the generator may produce an incorrect response, even though the optimal document is still present within the top-N set.

From an optimization perspective, this phenomenon reveals a deeper mismatch between reranking and generation. On the one hand, the

generator is highly sensitive to fine-grained ranking results; on the other hand, learning an exact total order over multiple highly relevant documents is inherently harder than learning a relaxed ordering for the reranker. For instance, precisely ranking the top three relevant documents as positions 1, 2, and 3 is substantially more challenging than merely ensuring that these documents appear within the top few positions in any order. This discrepancy between the difficulty of ranking optimization and the generator’s sensitivity to reranking results poses significant challenges for the stability and generalization of RAG frameworks in practice.

To address these issues, we reformulate RAG as a cooperative multi-agent decision-making problem, called Cooperative Retrieval-Augmented Generation (CoRAG) (Figure 1(b)). Unlike ranking-centric RAG frameworks that enforce an asymmetric dependency from reranking to generation, CoRAG jointly optimizes the reranker and the generator under a shared, task-oriented reward. This cooperative formulation relaxes the generator’s asymmetric dependency on reranking quality, encouraging the reranker to learn an accurate total order over highly relevant documents and simultaneously training the generator to robustly utilize information rather than relying on a single overly strict ranking. As a result, CoRAG mitigates the generator’s asymmetric dependency to fine-grained reranking results and improves generation stability. Experimental results show that, despite being trained on only around 10K PopQA samples, CoRAG significantly outperforms baselines and generalizes well across multiple datasets and tasks.

In summary, this paper makes the following contributions:

- We model RAG as a cooperative decision-making problem, treating the reranker and generator as peer agents optimized toward a shared objective.
- We propose a joint optimization scheme that mitigates asymmetric dependency of the generator on ranking results.
- Extensive experiments demonstrate improved robustness and generalization of our CoRAG.

## 2 Problem Definition

A typical RAG system consists of two components: a retriever and a generator (Lewis et al., 2020; Oche et al.). Given an input query, the retriever retrieves a candidate document set from a

large external corpus, which is further refined by a reranker before being used by the generator to generate the final response. In this work, we focus on the reranker in the retrieval stage together with the generator, as reranking plays a crucial role in improving the relevance and faithfulness of generation (Sharma, 2025). We model RAG as a cooperative multi-agent decision-making problem, treating the reranker and the generator as two cooperative agents that jointly determine the final generation outcome, and train them with the goal of maximizing a shared task-oriented objective defined on the generated response.

Specifically, given a query  $q$  and a candidate document set  $\mathcal{D}$ , the reranker  $\mathcal{S}_\theta$  reranks and selects a set of documents  $D \subseteq \mathcal{D}$ , and the generator  $\mathcal{G}_\phi$  generate an response  $\hat{a}$  conditioned on  $q$  and  $D$ . The learning objective is defined as

$$\max_{\theta, \phi} \mathbb{E}_{D \sim \mathcal{S}_\theta(\cdot|q, \mathcal{D}), \hat{a} \sim \mathcal{G}_\phi(\cdot|q, D)} [R(a^*, \hat{a})], \quad (1)$$

where  $\theta$  and  $\phi$  are the parameters of the reranker and generator respectively.  $R(a^*, \hat{a})$  denotes a task-oriented reward defined on the generated response. This formulation emphasizes that the reranker is optimized not for ranking accuracy, but for its contribution to downstream generation quality. Similarly, the generator is trained to produce outputs that effectively utilize the provided document configuration to maximize the same task-oriented reward. By aligning both components to the same outcome-oriented objective, the reranker and generator are encouraged to cooperate, ensuring that document reranking and generation in concert to improve the final response.

## 3 Method

An overview of our proposed Cooperative Retrieval-Augmented Generation (CoRAG) is shown in Figure 2. In the following sections, we first describe the reranker and generator in our CoRAG, and then discuss their joint optimization.

### 3.1 The Reranker

The reranker aims to refine the retrieved candidate document set  $\mathcal{D}$  by re-ranking the documents according to their relevance, and selects a subset  $D$  to provide to the downstream generator for response generation. Specifically, we obtain  $D$  in the following steps:

Given a query  $q$  and a document  $d_i$  in the candidate document set  $\mathcal{D} = \{d_1, d_2, \dots, d_N\}$ , we

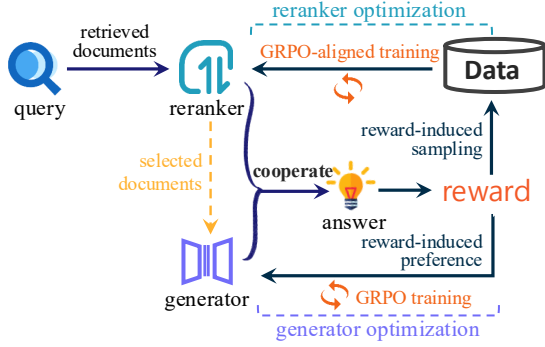


Figure 2: CoRAG overview. The reranker and generator cooperate to generate responses. The task-oriented reward derived from response guides GRPO-aligned training of the reranker and GRPO optimization of the generator.

compute the relevance score of document  $d_i$  to  $q$  by

$$s_i = \mathcal{S}_\theta(q, d_i), \quad (2)$$

$\mathcal{S}_\theta$  denotes the reranker and we implemented with BGE-Reranker (Multi-Granularity, 2024).

We select a subset of top- $K$  documents according to the reranker scores and feed them into the generator.

$$D = \{d_i \mid i \in \text{Top-K}(\{s_1, \dots, s_N\})\} \quad (3)$$

### 3.2 The Generator

The Generator is responsible for generating the response  $\hat{a}$  based on the selected document  $D$  and the query  $q$ . Specifically, we obtain the predicted response  $\hat{a}$  with the following steps:

Given a query  $q$  and the top- $K$  documents  $D$ , we input them into the generator model  $\mathcal{G}_\phi$  to generate the final response:

$$\hat{a} = \mathcal{G}_\phi(q, D), \quad (4)$$

where  $\mathcal{G}_\phi$  is the generator, typically implemented as an autoregressive language model (Radford et al., 2018; Brown et al., 2020) that generates the response token by token conditioned on both the query  $q$  and the document  $D$ .

### 3.3 The Optimization

The reranker and the generator are optimized under a shared task-oriented reward. From a multi-agent perspective, we treat the reranker and the generator as two cooperative agents with distinct decision roles. The reranker determines which documents to attend to, while the generator decides how to

synthesize the final response based on the selected documents. The shared reward, defined on the final response, couples their behaviors and enables coordinated optimization toward task-oriented success.

$$r = R(a^*, \hat{a}), \quad (5)$$

where  $R(a^*, \hat{a})$  denotes a task-oriented evaluation function, which we implement to return 1 if the generated response  $\hat{a}$  contains the ground-truth response  $a^*$ , and 0 otherwise. In the following, we detail the optimization of the reranker and the generator under this shared task-oriented reward.

**Reranker Optimization.** Unlike standard learning-to-rank settings (Casalegno, 2022), document-level supervision is not directly available. To address that, we transform task-oriented rewards into document-level stochastic preference signals for reranker optimization. Since these signals indicate how documents collectively contribute to task success, we optimize the reranker in a group-relative preference aligned with GRPO (Shao et al., 2024).

First, to transform delayed task-level rewards into document-level stochastic preference signal, for each document  $d_i$  in the top- $K$  document set  $D$  at training iteration  $t$ , we assign a binary success signal  $l^{(t)}(q, d_i) \in \{0, 1\}$ , indicating whether the generated response achieved task success when  $d_i$  was included. Although this signal provides only a coarse and noisy approximation of individual document contribution under multi-document conditioning, it enables scalable credit assignment without requiring explicit supervision.

Based on these signals, we estimate the expected task success associated with each document  $d_i$ :

$$\bar{l}(q, d_i) = \frac{1}{T} \sum_{t=1}^T l^{(t)}(q, d_i), \quad (6)$$

where  $T$  denotes the number of historical iterations. To account for uncertainty, we map  $\bar{l}(q, d_i)$  to a smoothed Bernoulli parameter

$$p(q, d_i) = \alpha + (1 - 2\alpha) \cdot \bar{l}(q, d_i), \quad (7)$$

where  $\alpha \in (0, 0.5)$ , and we sample a stochastic preference label by:

$$p_i \sim \text{Bernoulli}(p(q, d_i)), \quad p_i \in \{0, 1\}. \quad (8)$$

This stochastic preference label serves as a surrogate feedback signal for reranker optimization,

enabling us to optimize the reranker in a group-  
relative preference framework aligned with GRPO.  
To conduct the optimization, we model the reranker  
as a deterministic scoring function  $\mathcal{S}_\theta(q, d_i)$  that  
induces a distribution over candidate documents:  
 $\pi_\theta(d_i | q, \mathcal{D}) \propto \exp(\mathcal{S}_\theta(q, d_i))$ . The group-  
relative advantage  $\hat{A}(q, d)$  is derived from these  
preference labels: positive labels ( $p_i = 1$ ) indi-  
cate relative advantage for document  $d_i$  within the  
group, while a negative label indicates a disadvan-  
tage. This leads to the GRPO objective:

$$\mathcal{L}_{\text{GRPO}}^r = -\mathbb{E}_{d_i \sim \pi_\theta} \left[ \hat{A}(q, d_i) \log \pi_\theta(d_i | q, \mathcal{D}) \right], \quad (9)$$

However, directly optimizing the stochastic ob-  
jective suffers from high variance under noisy  
credit estimates. We therefore reinterpret the group-  
relative advantages as inducing pairwise prefer-  
ences among candidate documents: documents  
with higher empirical success rates (more likely  
to receive positive labels) should be ranked higher.  
This allows us to reduce GRPO to a deterministic  
learning-to-rank problem. Constructing positive  
and negative sets  $\mathcal{D}^+$  and  $\mathcal{D}^-$  from the sampled  
labels, we adopt a margin-based pairwise ranking  
surrogate:

$$\mathcal{L}_{\text{rank}} = \sum_{d_i^+ \in \mathcal{D}^+} \sum_{d_j^- \in \mathcal{D}^-} \max(0, s_\theta(q, d_j^-) - s_\theta(q, d_i^+) + \gamma), \quad (10)$$

where  $\gamma$  is the margin hyperparameter. Although  
this reduction is not a strict equivalence to GRPO,  
the ranking loss preserves the group-relative pref-  
erences induced by the stochastic labeling scheme  
and the underlying GRPO objective in expectation.  
It thus provides a stable and efficient surrogate for  
optimizing the reranker while maintaining align-  
ment with task-oriented rewards.

**Generator Optimization.** The generator defines  
a conditional generation policy  $\pi_\phi(\hat{a} | q, D)$ ,  
where  $D$  is the set of top-K documents selected  
by the reranker. The generator is optimized using  
standard GRPO for conditional text generation:

$$\mathcal{L}_{\text{gen}} = -\mathbb{E}_{\hat{a} \sim \pi_\phi} \left[ \hat{A}(\hat{a}) \log \pi_\phi(\hat{a} | q, D) \right], \quad (11)$$

where the group-relative advantage  $\hat{A}(\hat{a})$  is com-  
puted from the task-oriented reward  $r$  using a base-  
line that compares the current response to other  
responses in the same training batch.

Overall, both the reranker and the generator are  
optimized with respect to the same task-oriented

reward, enabling them to cooperatively improve  
retrieval relevance and generation quality.

## 4 Experiments

We conduct extensive experiments to evaluate the  
performance of CoRAG, and we particularly focus  
on the research questions: (i) Is CoRAG more ef-  
fective than existing methods (**RQ1**)? (ii) How im-  
portant are the reranker and the generator to overall  
performance, and is it necessary to jointly optimize  
them (**RQ2**)? (iii) How does performance vary  
with the number of documents used (**RQ3**)? (iv)  
Does CoRAG generalize to other tasks (**RQ4**)? (v)  
Does the CoRAG generator produce high-quality  
outputs as judged by other LLMs (**RQ5**)?

### 4.1 Experimental Setting

**Datasets** We evaluate our method on multiple  
knowledge-intensive benchmarks, following the  
dataset setup of recent works (Wei et al., 2024).  
Specifically, we evaluate on PopQA (Mallen et al.,  
2023), TriviaQA (Joshi et al., 2017), Natural Ques-  
tions (NQ) (Kwiatkowski et al., 2019), ASQA (Stel-  
makh et al., 2022). Additionally, we include 2Wiki-  
MultiHopQA (Ho et al., 2020), a Wikipedia-based  
cross-document multi-hop QA dataset that requires  
models to reason across multiple documents to de-  
rive responses.

Table 1: Dataset statistics.

Dataset	Train	Test
PopQA	12,868	1,399
TriviaQA	78,785	11,313
NaturalQuestions	79,168	3,610
ASQA	4,353	948
2WikiMultiHopQA	167,454	12,576

**Evaluation Metrics.** Following InstructRAG  
(Wei et al., 2024), We report correctness, citation  
precision, and citation recall (Gao et al., 2023a)  
for ASQA. For the other datasets, we use accuracy  
measures whether the ground-truth responses are  
included in the generated outputs.

**Baselines.** Following InstructRAG (Wei et al.,  
2024), We compare our CoRAG with three  
groups of method: Baselines without Retrieval  
(relying solely on parametric knowledge, includ-  
ing ChatGPT, Llama-3-Instruct<sub>8B</sub>, and Llama-3-  
Instruct<sub>70B</sub>), RAG without Training (leveraging  
retrieved documents via in-context learning or  
prompting, such as In-Context RALM (Ram et al.,

Table 2: Overall results of our method and baselines on five benchmarks. Baseline results are reported in InstructRAG (Wei et al., 2024). - indicates the results are not reported not applicable. The best and second-best performances are highlighted in bold and with an underline, respectively. Notably, our models are trained only on PopQA (Mallen et al., 2023), while all other datasets are used exclusively for evaluation.

Method	PopQA (acc)	TriviaQA (acc)	NQ (acc)	2WikiMultiHopQA (acc)	ASQA		
					(em)	(pre)	(rec)
<i>Baselines w/o Retrieval</i>							
<b>Vanilla Zero-shot Prompting</b>							
ChatGPT	29.3	74.3	–	–	35.3	–	–
Llama-3-Instruct <sub>8B</sub>	22.8	69.4	46.6	45.6	30.6	–	–
Llama-3-Instruct <sub>70B</sub>	28.9	80.6	57.9	57.5	39.1	–	–
<i>RAG w/o Training</i>							
<b>In-Context RALM (Ram et al., 2023)</b>							
ChatGPT	50.8	65.7	–	–	40.7	65.1	76.6
Llama-3-Instruct <sub>8B</sub>	62.3	71.4	56.8	43.4	40.0	62.1	66.4
Llama-3-Instruct <sub>70B</sub>	63.8	76.3	60.2	51.2	43.1	62.9	67.6
<b>Few-Shot Demo. w/ Instruction</b>							
Llama-3-Instruct <sub>8B</sub>	63.1	74.2	60.1	45.3	42.6	55.0	64.4
Llama-3-Instruct <sub>70B</sub>	63.9	79.1	62.9	53.9	45.4	49.3	57.1
<b>InstructRAG-ICL (Wei et al., 2024)</b>							
Llama-3-Instruct <sub>8B</sub>	64.2	76.8	62.1	50.4	44.7	70.9	74.1
Llama-3-Instruct <sub>70B</sub>	65.5	81.2	66.5	57.3	47.8	69.1	71.2
<i>RAG w/ Training</i>							
<b>Vanilla Supervised Fine-tuning</b>							
Llama-3-Instruct <sub>8B</sub>	61.0	73.9	56.6	56.1	43.8	–	–
<b>Self-RAG (Asai et al., 2024)</b>							
Llama-2 <sub>7B</sub>	55.8	68.9	42.4	35.9	30.0	66.9	67.8
Llama-2 <sub>13B</sub>	56.3	70.4	46.4	36.0	31.4	<b>70.3</b>	<b>71.3</b>
Llama-3-Instruct <sub>8B</sub>	55.8	71.4	42.8	32.9	36.9	<u>69.7</u>	69.7
<b>RetRobust (Yoran et al., 2023)</b>							
Llama-2 <sub>13B</sub>	–	–	39.6	51.5	–	–	–
Llama-3-Instruct <sub>8B</sub>	56.5	71.5	54.2	54.7	40.5	–	–
<b>InstructRAG-FT (Wei et al., 2024)</b>							
Llama-3-Instruct <sub>8B</sub>	<u>66.2</u>	<u>78.5</u>	<u>65.7</u>	<u>57.2</u>	<b>47.6</b>	65.7	<u>70.5</u>
<b>CoRAG</b>							
Llama-3-Instruct <sub>8B</sub>	<b>71.2</b>	<b>81.0</b>	<b>72.4</b>	<b>58.2</b>	<u>45.8</u>	54.9	48.9

2023), Few-shot Demonstration with Instruction, and InstructRAG-ICL (Wei et al., 2024)), and RAG with Training (involving explicit training in the retrieval-augmented framework, including Self-RAG (Asai et al., 2024) for iterative improvement via self-retrieval, RetRobust (Yoran et al., 2023) for noise robustness training, and InstructRAG-FT (Wei et al., 2024) for instruction-following fine-tuning with retrieval).

**Implementation Details** We adopt BGE-reranker-v2-m3 (Multi-Granularity, 2024) as the reranker and Llama-3-Instruct-8B (Dubey et al., 2024) as the generator. During training, we use LLaMA3 to provide coarse annotations for positive and negative documents in the training dataset, which helps the reranker learn more effectively by alleviating the sparsity of stochastic preference labels at the early stages of training. We set the learning rate of reranker to  $5e-5$ , and the learning rate of generator to  $1e-5$ . Both reranker and generator adopt the LoRA fine-tuning strategy for efficient parameter updating. we set  $\gamma = 1$  in

Eq.(10), and the top-K selected documents  $K = 1$  to precisely attribute the impact of individual documents on task success. During inference, we set  $K = 3$  for PopQA, TriviaQA, NQ and ASQA, and  $K = 7$  for 2WikiMultiHopQA. The generator uses a temperature coefficient of 0.7 to balance generation diversity and stability. It is worth noting that our CoRAG are trained only on PopQA (Mallen et al., 2023), while all other datasets are used solely for evaluation, both due to our limited computational resources and to enable assessment of the generalization ability of our CoRAG.

## 4.2 Main Results (RQ1)

The results are reported in Table 2. As it is shown, our method (CoRAG) achieves outstanding performance. Notably, unlike the InstructRAG series of methods that perform separate training on each dataset, our model is trained only on the PopQA. Even so, our method yields state-of-the-art results on four core tasks (PopQA, TriviaQA, NQ and 2WikiMultiHopQA), with the correspond-

ing accuracy rates reaching 71.2%, 81.0%, 72.4% and 58.2% respectively, which significantly outperform existing methods such as RetRobust and InstructRAG-FT. This can be attributed to the joint optimization framework, in which the reranker progressively selects more relevant documents and the generator continuously improves its ability to extract and integrate information from them.

However, our CoRAG underperforms on the ASQA dataset. We attribute this primarily to the task discrepancy: ASQA requires synthesizing answers from multiple sources and handling ambiguous questions, which differs substantially from the factoid-style, single-answer questions prevalent in our training data (PopQA). Consequently, the retrieval-generator synergy optimized on factoid QA does not generalize effectively to this more complex, multi-answer setting.

### 4.3 Ablation Study (RQ2)

To investigate the importance of the reranker and the generator to overall performance, as well as whether it is necessary to jointly optimize them, we conduct ablation studies. Specifically, we have four variants:

- **Rtrain**: only finetune the reranker with the labels annotated by Llama-3-Instruct-8B.
- **Gtrain**: only finetune the generator with GRPO.
- **RGReplace**: Replace the generator of CoRAG with Llama-3-Instruct-8B, and replace the reranker with BGE-Reranker in inference.
- **GReplace**: Replace the the generator of CoRAG with Llama-3-Instruct-8B in inference.

The result presented in Table 3, we can observe that: (1) The individually trained RTrain and GTrain underperform CoRAG across both the PopQA and NQ datasets, which implies that a multi-agent cooperative optimization formulation may better capture the interaction between the reranker and the generator, leading to improved performance compared to independent optimization. (2) Comparing RGReplace, GReplace and CoRAG reveals that the generator contributes more significantly to performance improvement, whereas the reranker plays a relatively limited role. We attribute this phenomenon to the joint optimization. Because the reranker and generator are jointly optimized, the generator may achieve strong performance even under suboptimal reranking signals,

Table 3: Ablation study.

Dataset	PopQA	NQ
RTrain	66.19	62.71
GTrain	66.54	63.37
RGReplace	65.83	63.21
GReplace	66.26	62.46
CoRAG	<b>71.26</b>	<b>72.49</b>

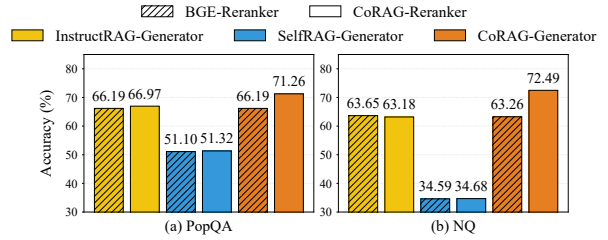


Figure 3: Cross-validation results with different reranker and generator combinations (Top-3 setting).

potentially weakening the learning pressure on the reranker and reflecting a trade-off in the current design.

To examine this hypothesis, we conduct cross-component experiments by swapping rerankers and generators between CoRAG and other RAG frameworks (Self-RAG and InstructRAG). As shown in Figure 3, replacing the reranker with our reranker while keeping other generators (self-RAG and InstructRAG) fixed yields only marginal improvements and occasionally leads to performance degradation. This suggests that the standalone effectiveness of our reranker is relatively limited. However, when paired with our generator, CoRAG consistently outperforms the variant that combines our generator with BGE-reranker. This contrast reveals an inherent tension between reranking effectiveness and generator sensitivity: jointly optimizing the reranker and generator encourages the generator to rely less on fine-grained ranking signals, which in turn limits the observable gains from further improving the reranker in isolation. Nevertheless, the strong performance of CoRAG indicates that its reranker and generator are well aligned and mutually reinforcing when optimized together.

### 4.4 Top-N Analysis (RQ3)

To analyze how the number of documents provided to the generator affects performance, we evaluate top-1(T1), top-3(T3), and top-5(T5) settings across different datasets. The results are presented in Figure 4, which reveal three insights: (1) CoRAG outperforms InstructRAG and RetRobust under all document count settings on both datasets, this indi-

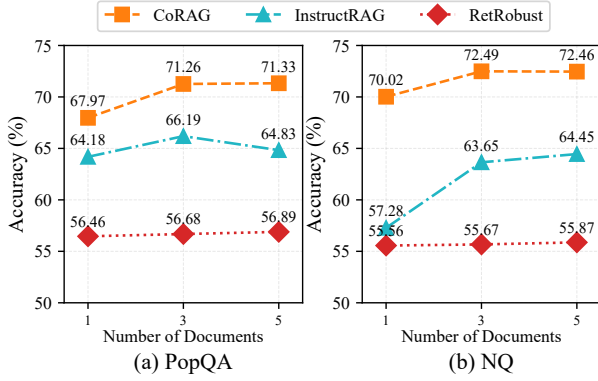


Figure 4: Impact of the document number.

Table 4: The cross-task evaluation. The InstructRAG in the table is trained on the Pop dataset. WTQ represents WikiTable Questions. The best and second-best performances are highlighted in bold and with an underline.

Metric	Llama3	InstructRAG	CoRAG
<b>human-eval</b>			
pass@1	<b>64.45</b>	55.85	<u>63.96</u>
pass@10	<u>83.53</u>	76.82	<b>85.97</b>
<b>human-eval+</b>			
pass@1	<u>56.58</u>	49.26	<b>57.43</b>
pass@10	<u>77.43</u>	70.12	<b>79.26</b>
<b>WTQ</b>			
accuracy	49.10	<b>68.57</b>	<u>68.11</u>

449 cates that CoRAG can effectively utilize the effective  
 450 information of documents under different docu-  
 451 ment number; (2) on PopQA, InstructRAG’s per-  
 452 formance declines as document number increases  
 453 (e.g., dropping from 66.19% at T3 to 64.83% at T5),  
 454 implying more documents may trigger the noise  
 455 of less relevant documents, and leading to hallu-  
 456 cinations. In contrast, our CoRAG demonstrates  
 457 a consistent upward trend in performance on both  
 458 PopQA and NQ, showing strong robustness.

#### 4.5 Cross-Task Evaluation (RQ4)

460 To further explore the cross-domain generalization  
 461 capability of CoRAG, we conduct additional ex-  
 462 periments on other tasks, following (Wei et al.,  
 463 2024). Specifically, we evaluate the generator of  
 464 our CoRAG and other generator (Llama3 and In-  
 465 structRAG) on code generation datasets HumanEval,  
 466 HumanEval+ (Chen, 2021), and table question  
 467 answering dataset WikiTable Questions (Pasupat and  
 468 Liang, 2015). Results are summarized in Table 4.

469 As shown in the table Table 4, compared with the  
 470 Llama3 and InstructRAG, our generator achieves  
 471 the best or second-best performance across all  
 472 tasks: on HumanEval and HumanEval+, CoRAG  
 473 achieves the strongest results on pass@10 and

Table 5: Evaluation with LLM-as-a-judge.

Method	InstructRAG	RetRobust	CoRAG
Pattern-based	66.19	56.68	<b>71.26</b>
Llama	75.63	25.52	<b>89.21</b>
GPT	60.83	54.90	<b>65.90</b>
DeepSeek	59.69	34.10	<b>64.26</b>
Qwen	60.69	57.26	<b>66.12</b>
Average	64.60	45.69	<b>71.35</b>

474 achieves the best pass@1 performance on the more  
 475 challenging HumanEval+, indicating improved  
 476 generation quality and robustness on code genera-  
 477 tion. On WTQ, a table-based question answering  
 478 benchmark, the generator of CoRAG achieves com-  
 479 petitive accuracy, closely matching the strongest  
 480 baseline while substantially outperforming Llama3.  
 481 Overall, these results demonstrate that the genera-  
 482 tor of CoRAG remains effective across both code  
 483 generation and table-based reasoning tasks.

#### 4.6 Evaluation with LLM-as-a-judge (RQ5)

484 Following InstructRAG (Wei et al., 2024), We use  
 485 LLMs as a judgment to further evaluate the genera-  
 486 tion of our CoRAG. Specifically, for a given ques-  
 487 tion and documents ranked by reranker, responses  
 488 are generated by generator and its quality is as-  
 489 sessed using different LLMs. The specific prompt  
 490 template is shown in Appendix B. We choose  
 491 Llama, GPT, DeepSeek, Qwen as the judge. For  
 492 LLaMa, we use version Llama-3.1-8b-instruction.  
 493 For GPT, we use version gpt-4o. For DeepSeek,  
 494 we use version deepseek-v3.2. For Qwen, we use  
 495 version qwen3-vl-235b-a22b-instruction. The re-  
 496 sults are presented in Table 5. As shown in Table  
 497 5, CoRAG consistently outperforms both Instruc-  
 498 tRAG and RetRobust across all evaluation settings,  
 499 regardless of the LLM used as the judge. CoRAG  
 500 achieves the highest scores under every evaluator,  
 501 including pattern-based metrics as well as Llama,  
 502 GPT, DeepSeek, and Qwen, indicating robust and  
 503 consistently preferred generation quality. Notably,  
 504 the performance gap is most pronounced under the  
 505 Llama judge, where CoRAG receives higher scores  
 506 than other LLMs. We attribute this effect to the  
 507 fact that CoRAG is fine-tuned from Llama, which  
 508 may lead to a higher alignment between CoRAG’s  
 509 outputs and the Llama-based judge.  
 510

## 5 Related Work

511 RAG enhances LLMs’ performance by fusing ex-  
 512 ternal documents and is the mainstream solution  
 513 for knowledge-intensive tasks. Existing research  
 514

could be categorized into three types based on core efficiency-enhancing mechanisms (Gao et al., 2023b; Zhao et al., 2023). The first category is data-driven methods, focusing on the mining and reconstruction of information at the query or document level: Decomposition Prompting (Khot et al., 2022), which splits complex tasks through prompt engineering. EviNoteRAG (Dai et al., 2025) annotates uncertain information in documents with a note-taking-first approach. HtmlRAG (Tan et al., 2025) uses HTML instead of plain text to preserve semantic structural information. The second category is model-driven methods, which improve a model’s ability to interpret, filter and utilize retrieved documents through fine-tuning. Representative works include: RetRobust (Yoran et al., 2023) enhances robustness through contrastive training with positive and negative samples. InstructRAG (Wei et al., 2024) explicitly learns the denoising process through self synthesized reasoning. DynamicRAG (Sun et al., 2025) optimizes the order and quantity of documents used to train the reranker through generator feedback. The third category is strategy-driven methods, also called agentic RAG, which introduces agentic behaviors to dynamically adjust retrieval and generation strategies. FLARE (Jiang et al., 2023) triggers lookahead retrieval when encountering uncertain tokens during generation. SelfRAG (Asai et al., 2024) introduces a reflection module to synchronously and dynamically adjust both retrieval and generation. MA-RAG (Nguyen et al., 2025) decomposes workflows into sub-agents via CoT. ComposeRAG (Wu et al., 2025) enables modular agent composition.

From a high-level perspective, CoRAG could be considered similar to agentic RAG, as it models the reranker and generator as cooperative agents. However, unlike many existing agentic RAG methods that treat the retriever (reranker) and generator as separate, modular agents often requiring separate training or intricate prompting for collaboration, CoRAG frames them as a unified policy, eliminating the need for handcrafted coordination mechanisms and allows both modules to adaptively specialize in a task-aware manner.

## 6 Conclusion and Limitations

In this paper, we reformulate Retrieval-Augmented Generation (RAG) as a cooperative multi-agent decision-making problem and propose Cooperative Retrieval-Augmented Generation (CoRAG). Un-

like conventional RAG frameworks, where the generator exhibits an asymmetric dependency on the reranker and generation quality is highly sensitive to the reranking results, CoRAG treats the reranker and generator as two peer decision-makers. Specifically, the reranker decides which documents, and in what organization, to present to the generator, while the generator determines how to effectively utilize the provided information to accomplish the generation task. Through this cooperative formulation, CoRAG alleviates the generator’s reliance on overly strict document ordering and enables more robust coordination between retrieval and generation, leading to improved overall performance.

## 7 Limitations

Although modeling RAG as a cooperative multi-agent framework improves robustness by reducing the generator’s sensitivity to reranking results, this design may also attenuate the impact of further reranker improvements on generation quality. This reflects an inherent tension between reranking effectiveness and generation sensitivity under joint optimization, which we leave for future exploration.

## References

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-rag: Learning to retrieve, generate, and critique through self-reflection.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Francesco Casalegno. 2022. Learning to rank: A complete guide to ranking using machine learning. *Towards Data Science*.
- Mark Chen. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Yuin Dai, Guoqing Wang, Yuan Wang, Kairan Dou, Kaichen Zhou, Zhanwei Zhang, Shuo Yang, Fei Tang, Jun Yin, Pengyu Zeng, and 1 others. 2025. Evinote-rag: Enhancing rag models via answer-supportive evidence notes. *arXiv preprint arXiv:2509.00877*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

614	Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen.	Multi-Linguality Multi-Functionality Multi-Granularity.	669
615	2023a. Enabling large language models to generate	2024. M3-embedding: Multi-linguality, multi-	670
616	text with citations. <i>arXiv preprint arXiv:2305.14627</i> .	functionality, multi-granularity text embeddings	671
		through self-knowledge distillation.	672
617	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang	Thang Nguyen, Peter Chin, and Yu-Wing Tai. 2025.	673
618	Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun,	Ma-rag: Multi-agent retrieval-augmented generation	674
619	Haofen Wang, and Haofen Wang. 2023b. Retrieval-	via collaborative chain-of-thought reasoning. <i>arXiv</i>	675
620	augmented generation for large language models: A	<i>preprint arXiv:2505.20096</i> .	676
621	survey. <i>arXiv preprint arXiv:2312.10997</i> , 2(1).		
622	Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara,	AJ Oche, AG Folashade, T Ghosal, and A Biswas. A	677
623	and Akiko Aizawa. 2020. Constructing a multi-hop	systematic review of key retrieval-augmented genera-	678
624	qa dataset for comprehensive evaluation of reasoning	tion (rag) systems: Progress, gaps, and future direc-	679
625	steps. <i>arXiv preprint arXiv:2011.01060</i> .	tions. arxiv 2024. <i>arXiv preprint arXiv:2409.15730</i> .	680
626	Pengyue Jia, Derong Xu, Xiaopeng Li, Zhaocheng Du,	Panupong Pasupat and Percy Liang. 2015. Composi-	681
627	Xiangyang Li, Yichao Wang, Yuhao Wang, Qidong	tional semantic parsing on semi-structured tables.	682
628	Liu, Maolin Wang, Huifeng Guo, and 1 others. 2025.	<i>arXiv preprint arXiv:1508.00305</i> .	683
629	Bridging relevance and reasoning: Rationale distil-	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya	684
630	lation in retrieval-augmented generation. In <i>Find-</i>	Sutskever, and 1 others. 2018. Improving language	685
631	<i>ings of the Association for Computational Linguistics:</i>	understanding by generative pre-training.	686
632	<i>ACL 2025</i> , pages 4242–4256.		
633	Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun,	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay,	687
634	Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie	Amnon Shashua, Kevin Leyton-Brown, and Yoav	688
635	Callan, and Graham Neubig. 2023. Active retrieval	Shoham. 2023. In-context retrieval-augmented lan-	689
636	augmented generation. In <i>Proceedings of the 2023</i>	guage models. <i>Transactions of the Association for</i>	690
637	<i>Conference on Empirical Methods in Natural Lan-</i>	<i>Computational Linguistics</i> , 11:1316–1331.	691
638	<i>guage Processing</i> , pages 7969–7992.		
639	Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	692
640	Zettlemoyer. 2017. Triviaqa: A large scale distantly	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	693
641	supervised challenge dataset for reading comprehen-	Zhang, YK Li, Yang Wu, and 1 others. 2024.	694
642	sion. <i>arXiv preprint arXiv:1705.03551</i> .	Deepseekmath: Pushing the limits of mathematical	695
643	Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao	reasoning in open language models. <i>arXiv preprint</i>	696
644	Fu, Kyle Richardson, Peter Clark, and Ashish Sab-	<i>arXiv:2402.03300</i> .	697
645	harwal. 2022. Decomposed prompting: A modular	Chaitanya Sharma. 2025. Retrieval-augmented genera-	698
646	approach for solving complex tasks. <i>arXiv preprint</i>	A comprehensive survey of architectures, en-	699
647	<i>arXiv:2210.02406</i> .	hancements, and robustness frontiers. <i>arXiv preprint</i>	700
648	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	<i>arXiv:2506.00054</i> .	701
649	field, Michael Collins, Ankur Parikh, Chris Alberti,	Weizhou Shen, Yeyun Gong, Yelong Shen, Song Wang,	702
650	Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-	Xiaojun Quan, Nan Duan, and Weizhu Chen. 2023.	703
651	ton Lee, and 1 others. 2019. Natural questions: a	Joint generator-ranker learning for natural language	704
652	benchmark for question answering research. <i>Trans-</i>	generation. In <i>Findings of the Association for Com-</i>	705
653	<i>actions of the Association for Computational Linguis-</i>	<i>putational Linguistics: ACL 2023</i> , pages 7681–7699.	706
654	<i>tics</i> , 7:453–466.		
655	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	Ivan Stelmakh, Yi Luan, Bhuwan Dhingra, and	707
656	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	Ming-Wei Chang. 2022. Asqa: Factoid ques-	708
657	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	tions meet long-form answers. <i>arXiv preprint</i>	709
658	täschel, and 1 others. 2020. Retrieval-augmented gen-	<i>arXiv:2204.06092</i> .	710
659	eration for knowledge-intensive nlp tasks. <i>Advances</i>	Jiashuo Sun, Xianrui Zhong, Sizhe Zhou, and Jiawei	711
660	<i>in neural information processing systems</i> , 33:9459–	Han. 2025. Dynamicrag: Leveraging outputs of large	712
661	9474.	language model as feedback for dynamic reranking	713
662	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,	in retrieval-augmented generation. <i>arXiv preprint</i>	714
663	Daniel Khashabi, and Hannaneh Hajishirzi. 2023.	<i>arXiv:2505.07233</i> .	715
664	When not to trust language models: Investigating	Jiejun Tan, Zhicheng Dou, Wen Wang, Mang Wang,	716
665	effectiveness of parametric and non-parametric mem-	Weipeng Chen, and Ji-Rong Wen. 2025. Htmrlrag:	717
666	ories. In <i>Proceedings of the 61st Annual Meeting of</i>	Htmrl is better than plain text for modeling retrieved	718
667	<i>the Association for Computational Linguistics (Vol-</i>	knowledge in rag systems. In <i>Proceedings of the</i>	719
668	<i>ume 1: Long Papers</i> ), pages 9802–9822.	<i>ACM on Web Conference 2025</i> , pages 1733–1746.	720

721 Zhepei Wei, Wei-Lin Chen, and Yu Meng. 2024. In-  
 722 structrag: Instructing retrieval-augmented genera-  
 723 tion via self-synthesized rationales. *arXiv preprint*  
 724 *arXiv:2406.13629*.

725 Ruofan Wu, Youngwon Lee, Fan Shu, Danmei Xu,  
 726 Seung-won Hwang, Zhewei Yao, Yuxiong He, and  
 727 Feng Yan. 2025. Composerag: A modular and com-  
 728 posable rag for corpus-grounded multi-hop question  
 729 answering. *arXiv preprint arXiv:2506.00232*.

730 Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan  
 731 Berant. 2023. Making retrieval-augmented language  
 732 models robust to irrelevant context. *arXiv preprint*  
 733 *arXiv:2310.01558*.

734 Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhen-  
 735 gren Wang, Yunteng Geng, Fangcheng Fu, Ling  
 736 Yang, Wentao Zhang, Jie Jiang, and Bin Cui. 2024.  
 737 Retrieval-augmented generation for ai-generated con-  
 738 tent: A survey. *arXiv preprint arXiv:2402.19473*.

739 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,  
 740 Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen  
 741 Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023.  
 742 A survey of large language models. *arXiv preprint*  
 743 *arXiv:2303.18223*, 1(2).

## 744 A The training of CoRAG

745 The training of CoRAG has been summarized as  
 746 Algorithm 1.

---

### Algorithm 1 Training of CoRAG

---

```

1: for  $t = 1$  to  $T$  do
2:   for query  $q \in Q$  do
3:     Compute score  $s_i = \mathcal{S}_\theta(q, d_i)$ ;
4:     Select top-K documents  $D$  according to scores;
5:     Generate answer:  $\hat{a} = \mathcal{G}_\phi(q, D)$ 
6:     Compute reward:  $R(a^*, \hat{a})$ ;
7:     Estimate the expected task success for  $d_i \in D$ ;
8:     Sample a stochastic preference label;
9:     Update reranker according to  $\mathcal{L}_{\text{rank}}$ ;
10:    Update generator according to  $\mathcal{L}_{\text{gen}}$ ;
11:   end for
12: end for

```

---

## 747 B Prompt Template

748 We provide the prompt used for LLM-as-a-Judge in  
 749 Table 6, and the prompt used by CoRAG in Table  
 750 7.

Table 6: Prompt of LLM-as-a-Judge

---

**Prompt:** You are an expert evaluator for large model responses. Your core task is to determine whether the large model’s response points to any of the correct answers.

Please conduct the evaluation based on the following information:

1. Question: {question}
2. List of correct answers: {answers}
3. Large model’s response: {response}

Evaluation Rules:

1. If the core information and key content of the large model’s response point to any answer in the correct answer list (semantic consistency is acceptable, no need for word-for-word matching), please output "yes".

2. If the large model’s response deviates from all correct answers, contains obvious errors, or fails to effectively respond to the question, please output "no".

3. You only need to output a single word: either "yes" or "no", without any additional redundant content.

---

**Output:** yes

---

Table 7: Prompt of Inference

---

**Prompt:** You are tasked with answering the given question by analyzing a set of documents. Please follow this STRICT TWO-STEP PROCESS:

---  
### STEP 1: Document Analysis

For each document:

- First, extract potentially relevant information from the original document. This includes facts, names, dates, or statements that may relate to the question, even if the connection is not immediately obvious.
- Then, explain the reason for the information extraction in your previous step based on the question. (e.g., how the document addresses the question's focus).

### STEP 2: Final Answer

- Summarize your answer to the question based on the analysis above. - If none of the documents are helpful or relevant, answer based on your own general knowledge. In that case, clearly state that you are doing so.

---  
Use the following format for your response:

### Step 1: Document Analysis

Document 1: - Extraction: ... - Explanation: ...

Document 2: - Extraction: ... - Explanation: ...

### Step 2: Final Answer

Well-supported answer, based on the relevant documents. If no relevant documents, answer based on general knowledge and say so explicitly.

---  
EXAMPLE: Question: Who is the author of The Mahdi?

### Step 1: Document Analysis

Document 1: - Extraction: 'Mahdi' is a thriller novel by Philip Nicholson written in 1981 under the identity of a. J. Quinnell ... - Explanation: This document directly states that the author of 'Mahdi' is Philip Nicholson. Upon re-examining the question "Who is the author of The Mahdi?", the document mentions the corresponding book title and provides the author information requested.

Document 2: - Extraction: 'The Mahdi' was published by Philip Nicholson in 1981 under the pen name A.J. Quinnell, establishing his presence in the thriller genre with a novel known for its gripping plot and enduring popularity... - Explanation: This document directly states that 'The Mahdi' was published by Philip Nicholson. Upon re-examining the question "Who is the author of The Mahdi?", the document mentions the corresponding book title and provides the author information requested.

### Step 2: Final Answer

Based on Documents 1 and 2, The answer to question 'Who is the author of The Mahdi?' is A.J. Quinnell, a pseudonym of Philip Nicholson.

---  
Now it is your turn to analyze the following documents and answer the given question by following the two-step process.

{context} {question}

---