

# Self-Critique Guided Iterative Reasoning for Multi-hop Question Answering

Zheng Chu<sup>1</sup>, Huiming Fan<sup>1</sup>, Jingchang Chen<sup>1</sup>, Qianyu Wang<sup>1</sup>, Mingda Yang<sup>1</sup>,  
Jiafeng Liang<sup>1</sup>, Zhongjie Wang<sup>1</sup>, Hao Li<sup>1</sup>, Guo Tang<sup>1</sup>, Ming Liu<sup>1,2\*</sup>, Bing Qin<sup>1,2</sup>

<sup>1</sup>Harbin Institute of Technology, Harbin, China

<sup>2</sup>Peng Cheng Laboratory, Shenzhen, China

{zchu,mliu}@ir.hit.edu.cn

## Abstract

Although large language models (LLMs) have demonstrated remarkable reasoning capabilities, they still face challenges in knowledge-intensive multi-hop reasoning. Recent work explores iterative retrieval to address complex problems. However, the lack of intermediate guidance often results in inaccurate retrieval and flawed intermediate reasoning, leading to incorrect reasoning. To address these, we propose Self-Critique Guided Iterative Reasoning (SiGIR), which uses self-critique feedback to guide the iterative reasoning process. Specifically, through end-to-end training, we enable the model to iteratively address complex problems via question decomposition. Additionally, the model is able to self-evaluate its intermediate reasoning steps. During iterative reasoning, the model engages in branching exploration and employs self-evaluation to guide the selection of promising reasoning trajectories. Extensive experiments on three multi-hop reasoning datasets demonstrate the effectiveness of our proposed method, surpassing the previous SOTA by 8.6%. Furthermore, our thorough analysis offers insights for future research. Our code, data, and models are available at Github<sup>1</sup>.

## 1 Introduction

With the advancement of large-scale pre-training techniques, LLMs have achieved remarkable performance on reasoning tasks with chain-of-thought (CoT) prompting (Wei et al., 2022; Chu et al., 2024b; OpenAI, 2023; DeepSeek-AI et al., 2024a). Moreover, when augmented with *Slow Thinking* mechanisms, several advanced models have demonstrated reasoning capabilities that either surpass or are comparable to those of human experts (DeepSeek-AI et al., 2024b; OpenAI, 2024).

Despite their remarkable capabilities, LLMs are prone to generating factually inaccurate re-

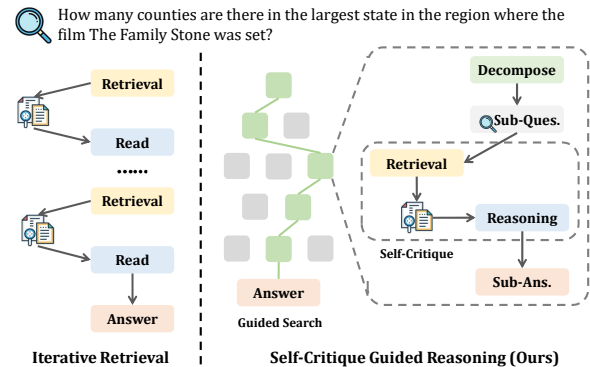


Figure 1: A comparison of iterative retrieval and self-critique guided reasoning (ours) in MHQA. Our method incorporates self-critique during iterative reasoning and conduct searches based on the received feedback.

sponses in knowledge-intensive tasks, primarily attributable to their limited access to domain-specific knowledge. This limitation becomes particularly pronounced in multi-hop question answering (MHQA) (Yang et al., 2018; Ho et al., 2020; Trivedi et al., 2022), which substantially compromises both the reasoning accuracy and factual reliability of the model’s outputs. Furthermore, the complexity of multi-hop questions presents additional challenges for global planning during the reasoning process

Retrieval-augmented generation is a promising strategy for mitigating the knowledge limitations of LLMs (Jiang et al., 2024). However, single-round retrieval is insufficient to encompass all relevant evidences to multi-hop questions (Lewis et al., 2020). To address this, iterative retrieval has been proposed (Trivedi et al., 2023; Asai et al., 2024). By alternating between retrieval and generation, it can acquire more comprehensive relevant knowledge, improving retrieval performance. Besides, some research focuses on the problems from the viewpoint of decomposition (Cao et al., 2023; Chu et al., 2024a; Su et al., 2024). It breaks down the complex multi-hop question into simpler sub-questions, thereby reducing the difficulty of reasoning and

\* Corresponding Author.

<sup>1</sup><https://github.com/zchuz/SiGIR-MHQA>

planning, and enabling more accurate retrieval.

Although the aforementioned methods have partially addressed some aspects of challenges, they still suffer from certain limitations. First, decomposing a complex question into sub-questions at the beginning of reasoning is highly challenging, and any mistakes can cause subsequent reasoning to deviate from the correct path. Second, the iterative retrieval approach struggles to handle complex problem planning and articulate retrieval intents clearly, resulting in inaccurate retrieval and reasoning. Finally, the absence of guidance for intermediate steps makes these methods susceptible to errors during reasoning, resulting in cascading errors.

To address these, our research focuses on following key aspects: (i) *iterative question decomposition* facilitates complex question planning and accurate retrieval. (ii) *intermediate feedback guided reasoning* expands the reasoning space and reduces cascading errors. Based on this motivation, we introduce **Self-Critique Guided Iterative Reasoning (SiGIR)**. SiGIR iteratively solves complex questions through sub-question decomposition and selects promising reasoning paths guided by the self-critique feedback of intermediate reasoning steps.

Specifically, during **training phase**, we synthesize iterative reasoning rationales with intermediate supervision from advanced LLMs. Thereafter, the model is equipped with the ability to perform question decomposition, trigger retrieval, knowledge reasoning, self-evaluation, and question reduction by means of end-to-end training. During **inference phase**, the model conducts iteration-level (step-level) beam search guided by self-feedback, as illustrated in Figure 1. At each reasoning step, the model performs branch expansion for sub-question decomposition, retrieval, and reasoning. It then self-evaluates the current reasoning status and retain the top-k promising candidates based on self-critique rewards. The model continues this iterative reasoning process until the final answer is derived. Finally, the reasoning trajectory with the highest cumulative reward is selected as the final answer.

We conduct extensive experiments on three knowledge-intensive MHQA datasets: HotpotQA, 2WikiMQA, and MuSiQue (Yang et al., 2018; Ho et al., 2020; Trivedi et al., 2022). DeepSeek-V2.5, Mistral, LLaMA2, and Qwen2.5 are used during experiments (DeepSeek-AI et al., 2024a; Jiang et al., 2023a; Touvron et al., 2023; Yang et al., 2024). Our method demonstrates significant advantages, with an average improvement of 8.6% compared

to the SOTA method. Furthermore, the thorough analysis highlights the effectiveness of our method in retrieval, self-critique, search, and efficiency.

## 2 Related Work

### 2.1 Retrieval-Augmented Generation

Retrieval-augmented generation (RAG) is an effective approach for mitigating hallucinations in LLMs (Gao et al., 2023; Huang et al., 2023). Retrieval-then-read, a widely-used approach in RAG, consists of query-based retrieval followed by conditional generation (Lewis et al., 2020; Borgeaud et al., 2022). However, its effectiveness is critically reliant on the quality of retrieval. To address this, various methods have been proposed. Query rewriting improves retrieval quality by reformulating retrieval queries (Ma et al., 2023; Chan et al., 2024; Xin et al., 2024). Active retrieval dynamically adjusts retrieval timing and intent using adaptive mechanisms (Jiang et al., 2023b; Asai et al., 2024). Multi-round retrieval enables the acquisition of comprehensive information relevant to the problem (Trivedi et al., 2023; Shao et al., 2023).

Our proposed method employs multi-round active retrieval to acquire more precise and comprehensive content for complex multi-hop questions.

### 2.2 Multi-hop Question Answering

MHQA requires gathering multiple pieces of evidence and systematic multi-step reasoning to arrive at the final answer. Generally, single-round retrieval is inadequate to acquire comprehensive information (Lewis et al., 2020; Borgeaud et al., 2022). Some methods adopt question decomposition to enhance accurate retrieval and reasoning (Cao et al., 2023; Chu et al., 2024a; Su et al., 2024). However, decomposing complex questions from the very beginning is particularly challenging. To this end, iterative reasoning approach is introduced, which employs iterative multi-round reasoning to gradually tackle complex problems (Trivedi et al., 2023; Yu et al., 2024; Li et al., 2024).

In this work, we utilize iterative question decomposition and reduction to simplify the decomposition process and improve retrieval effectiveness.

### 2.3 Inference-Time Scaling

Recent studies have found that increasing computational resources at inference can improve model performance, termed inference-time scaling (Snell

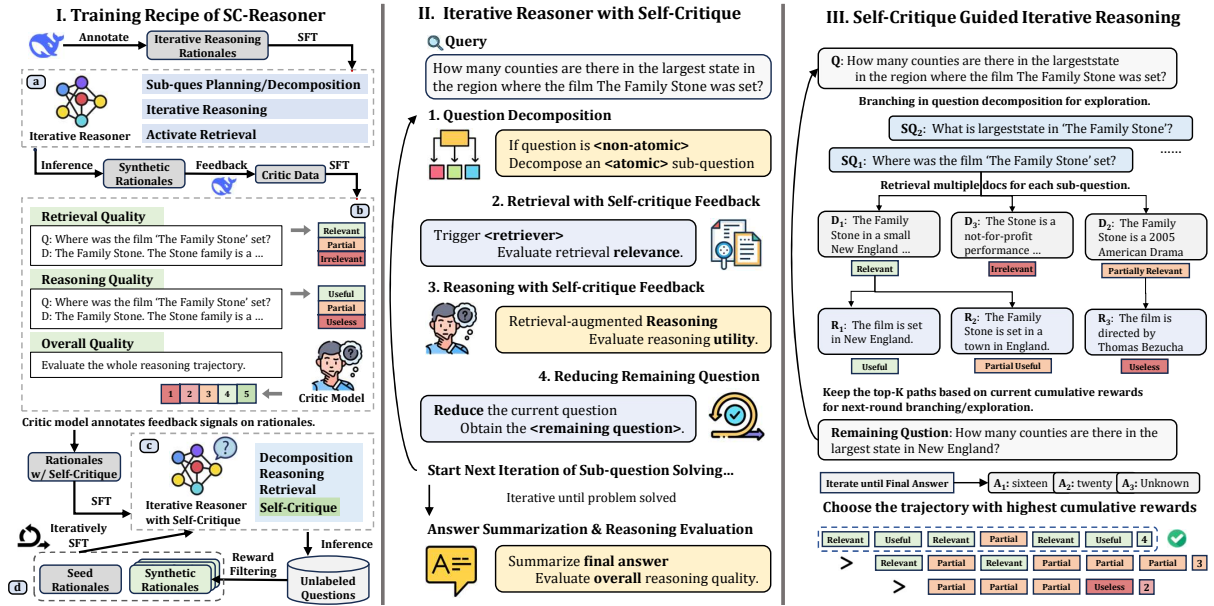


Figure 2: Overview of SiGIR. (I) Training recipe to get the iterative reasoner with self-critique,  $R_{sc}$ . (II) Reasoning characteristic of SC-Reasoner, which includes question decomposition, trigger retrieval, retrieval-augmented reasoning, self-evaluation, and question reduction. (III) SC-Reasoner engages in exploration during the iterative reasoning process and conducts reward-guided searches to select promising reasoning trajectories.

et al., 2024; Brown et al., 2024). Wang et al. (2023) employs a sampling ensemble technique to scale up, but the improvement is limited. Other methods adopt search strategies, such as Monte Carlo Tree Search, which yield significant improvements, but involve higher computational overhead (Zhang et al., 2024; Qi et al., 2024; Wang et al., 2024).

In contrast, we introduce reward-guided search that circumvents the overhead that simulations entail. It improves performance while ensuring reasonable costs, leading to efficient reasoning.

### 3 Self-Critique Guided Iterative Reasoning for Multi-hop QA

In this paper, we introduce Self-Critique Guided Iterative Reasoning (SiGIR) framework for multi-hop question answering, as illustrated in Figure 2. First, we synthesize iterative reasoning trajectories with intermediate supervision signals to train the model, which enables it to give self-critique of intermediate steps during reasoning. During reasoning, the model iteratively explores and selects the optimal trajectory guided by self-critique feedback.

In this section, we will outline the data synthesis pipeline and training recipe in Section 3.1, introduce the inference pattern of self-critique iterative reasoner in Section 3.2, and elaborate on the self-critique guided reasoning process in Section 3.3.

#### 3.1 Training Recipe

Our goal is to enable the model to perform iterative reasoning and to conduct self-critique of the intermediate reasoning steps. To achieve this, we first train an iterative reasoner  $R$ , followed by a critic model  $C$ . Subsequently, we use  $R$  and  $C$  to generate rationales and train a model  $R_{sc}$  capable of iterative reasoning and self-critique. Figure 2(I) and Algorithm 2 illustrate the training procedure.<sup>2</sup>

**Iterative Reasoner.** We use advanced LLMs to synthesize iterative reasoning rationales. Specifically, given a corpus  $X$ , where each instance  $x_i$  includes a question  $q_i$ , an answer  $a_i$ , and corresponding documents  $\{d_1, \dots, d_n\}$ . For each instance, we use few-shot prompts with the question and supporting documents as input to the LLM, which yields question decompositions and sub-question reasoning. Subsequently, we organize the sub-question decomposition, retrieval, and reasoning in an interleaved manner, and incorporate retrieval triggers before retrieved documents. We train the smaller model with the parsed training instances to get the vanilla iterative reasoner  $R$ .<sup>3</sup>

**Critic.** Due to the high cost of invoking LLMs for reward evaluation, we attempt to use a small model to clone the reward evaluation of LLMs. To

<sup>2</sup>Notations of training and inference are shown in Table 14.

<sup>3</sup>Prompts for data synthesis can be found in Appendix C.

this end, we use iterative reasoning trajectories to construct the data for critic training. Notably, to prevent the critic from overfitting to the generator’s training data, we select a small amount of non-overlapping corpus with  $X$  for critic data synthesis, denoted as  $X_{critic}$ . Additionally, we employ  $R$  instead of the LLM for reasoning synthesis, as  $R$  tends to make mistakes during reasoning. This helps ensure a balance between positive and negative samples. Subsequently, we employ the LLM to evaluate rewards on retrieval, reasoning, and overall quality, obtaining training data for critic model. We use it to train a separate critic model  $C$ .

**Self-Critique Reasoner.** We aim to incorporate the critique signals into the reasoning process. Besides, recent studies have found that the generator and critic can enhance each other, and merging them into a single model can reduce inference cost. To this end, we use the critic  $C$  to annotate the sub-processes of the reasoning trajectories. For retrieval evaluation, we take the query and retrieved documents as input, with retrieval relevance as outputs. For sub-question reasoning, the inputs include sub-questions, documents, and reasoning. For overall reward, we take the entire reasoning as input, and the output is the overall quality. Subsequently, we add special tokens following the sub-processes, enabling the model to generate self-critique signals throughout the reasoning process. We train the model with this corpus to obtain the reasoner with self-critique, denoted as SC-Reasoner, or  $R_{sc}$ .

SC-Reasoner is capable of reasoning and self-evaluation, facilitating data augmentation. Specifically, it can conduct reasoning on unlabeled questions and evaluate the quality of synthetic rationales, retaining high-quality rationales for subsequent training. This self-improvement pipeline can improve the performance in data scarcity scenarios. We use the standard causal language modeling objective for optimization during training.

### 3.2 Iterative Reasoner with Self-Critique

This section introduces the reasoning patterns of SC-Reasoner. Overall, SC-Reasoner is capable of question decomposition and planning, triggering retrieval, reasoning, and self-evaluating the reasoning process. As illustrated in Figure 2(II), given a question, the model first determines whether the question needs decomposition. If needed, the model decomposes one atomic sub-question at a time. Then, the model triggers external retrieval and conducts knowledge reasoning, along with self-evaluating

the quality of both retrieval and reasoning.

Whenever the model resolves a sub-question, it reduces the original question to identify the remaining unresolved part, thereby simplifying the overall question. The model iteratively repeats this process until the reduced atomic question is resolved, leading to the final answer. At this point, it conducts self-evaluation of the entire reasoning trajectory.

### 3.3 Self-Critique Guided Iterative Reasoning

After training, SC-Reasoner can conduct iterative reasoning and self-critique. Accordingly, we introduce self-critic guided iterative reasoning, an iteration-level beam search algorithm. This algorithm searches within the iterative reasoning process, employing self-rewards to guide the selection of promising reasoning candidates, as illustrated in Figure 2(III) and Algorithm 1. Generally, in each iteration, the model first expands the exploration space by branching based on the current reasoning prefix, and then selects top-k candidate trajectories guided by self-rewards. At the end of the reasoning, we use self-rewards to choose the final answer.

**Branch Exploration** Each iteration consists of the following steps: question decomposition, retrieval, reasoning, remaining question reduction, and self-critique. At timestamp  $t$ , the model branches out and expands downward through temperature sampling in question decomposition, retrieval, and reasoning, as shown in equations below.

$$q_{t+1} \sim p(y \mid q_t^{(r)}, x_{<t}) \quad (1)$$

$$s_{t+1} \sim p(y \mid q_{t+1}, d_i, x_{<t}) \quad (2)$$

where  $\sim$  is sampling,  $q_t^{(r)}$  is remaining question,  $q$  is sub-question,  $s$  is sub-question reasoning,  $p$  is probabilistic LM, and  $d$  is retrieved documents.

During the aforementioned process, the model also conducts self-evaluation and assigns rewards, as shown in Equ. (2,3). At the end of an iteration, the model will reduce the question for the next iteration based on current reasoning status  $x_{t+1}$ .

$$r_{t+1}^{\text{retr}} = p(y \mid q_{t+1}, d_i) \quad (3)$$

$$r_{t+1}^{\text{reas}} = p(y \mid q_{t+1}, d_i, s_{t+1}) \quad (4)$$

$$x_{t+1} = \oplus[x_t, q_{t+1}, d, r_{t+1}^{\text{retr}}, s_{t+1}, r_{t+1}^{\text{reas}}] \quad (5)$$

$$q_{t+1}^{(r)} = p(y \mid x_{t+1}) \quad (6)$$

**Candidate Selection** Through branch exploration, the model’s reasoning investigates a wide space. Subsequently, we aim to select the most



Methods	2WikiMQA					HotpotQA			MuSiQue			
	Overall	Brg.	Infer.	Comp.	B.C.	Overall	Brg.	Comp.	Overall	2hop	3hop	4hop
<i>Inference-based Methods</i>												
Chain-of-thought*	31.45	11.54	31.56	51.90	49.67	35.31	29.97	54.94	18.93	22.57	16.34	12.49
Self-Ask*	28.20	11.39	26.16	48.04	41.97	30.49	24.12	53.90	15.91	20.81	12.27	7.56
RAG	40.21	19.22	26.02	69.25	55.25	53.42	50.33	65.49	22.06	30.24	16.64	8.32
IRCoT	46.80	29.89	25.55	73.13	61.44	49.03	45.33	63.45	21.43	27.15	<u>18.73</u>	9.86
ProbTree <sup>†</sup>	59.90	-	-	-	-	50.40	-	-	27.00	-	-	-
BeamAggR <sup>†</sup>	63.20	-	-	-	-	55.20	-	-	<u>32.30</u>	-	-	-
<i>Training-based Methods</i>												
Self-RAG	31.60	16.91	14.06	50.39	49.00	54.60	52.51	62.75	22.00	30.04	13.21	<u>14.77</u>
RAG-Star <sup>†</sup>	46.80	-	-	-	-	<u>60.00</u>	-	-	30.70	-	-	-
GenGround <sup>†</sup>	50.21	-	-	-	-	52.26	-	-	27.36	-	-	-
Auto-RAG	56.06	44.12	40.94	79.44	59.48	50.03	46.96	62.02	22.08	29.32	15.90	11.54
DR-Distillation	<u>70.51</u>	<u>46.13</u>	<u>56.34</u>	<b>96.90</b>	<b>96.00</b>	58.06	<u>54.08</u>	<u>73.62</u>	22.74	<u>35.08</u>	12.06	6.57
<b>Ours</b>	<b>74.47</b>	<b>56.39</b>	<b>58.69</b>	<u>95.47</u>	<u>94.94</u>	<b>63.09</b>	<b>59.45</b>	<b>77.37</b>	<b>37.15</b>	<b>49.51</b>	<b>27.55</b>	<b>18.37</b>
- Qwen2.5	70.51	54.98	40.39	94.90	90.72	60.93	56.59	77.87	33.29	40.11	28.18	22.90
- LLaMA2	65.82	50.93	43.36	90.44	79.25	52.92	47.92	72.46	28.56	39.00	20.60	12.98

Table 1: Experimental results on three multi-hop reasoning datasets: 2WikiMQA, HotpotQA, and MuSiQue. Mistral-7B is used as the default backbone. Best and second results are highlighted by **bold** and underline. The evaluation metric is F1. \* indicates closed-book reasoning, and † denotes results are derived from the original paper.

promising candidates among them based on rewards. Selected candidates will serve as the starting point for the next iteration of reasoning. Specifically, we use the cumulative process reward as the criterion for candidate selection. The cumulative reward consists of reward from previous iteration and current iteration. We retain the candidates with top-k largest rewards, as shown in Equ. (7, 8). Once the reasoning is finished, we need to select one candidate as the final answer, using either cumulative reward of entire trajectory or outcome reward.

$$r_{t+1}^c = \underbrace{r_t^c}_{\text{previous reward}} + \underbrace{r_{t+1}^{\text{retr}} + r_{t+1}^{\text{reas}}}_{\text{current reward}} \quad (7)$$

$$\mathbf{x}_{t+1}^* = \text{top-k}(\mathbf{x}_{t+1}, \mathbf{r}_{t+1}^c) \quad (8)$$

$$r^c = \sum_{i=1}^n (r_i^{\text{retr}} + r_i^{\text{reas}}) \quad (9)$$

$$r^o = p(y | q^{(r)}, x_{<n}, a) \quad (10)$$

where  $r^c$  and  $r^o$  denote cumulative process reward and outcome reward, respectively.

## 4 Experimental Setup

### 4.1 Benchmarks

We evaluate SiGIR on three knowledge-intensive multi-hop reasoning datasets: HotpotQA (Yang et al., 2018), 2WikiMQA (Ho et al., 2020), and MuSiQue (Trivedi et al., 2022). Following previous work (Trivedi et al., 2023; Cao et al., 2023), we

use the development set as the test data. The experiments are conducted in the open-domain setting. We use token-level F1 as the evaluation metric.

### 4.2 Baselines

We select extensive baselines for comparison, including the following methods: Fewshot CoT (Wei et al., 2022), Self-Ask (Press et al., 2023), IR-CoT (Trivedi et al., 2023), ProbTree (Cao et al., 2023), BeamAggR (Chu et al., 2024a), RAG (Asai et al., 2024), RAG-Star (Jiang et al., 2024), GenGround (Shi et al., 2024), Auto-RAG (Yu et al., 2024), and DR-Distillation (Li et al., 2024).<sup>4</sup>

### 4.3 Implementation Details

We use DeepSeek-v2.5 (DeepSeek-AI et al., 2024a) as the advanced LLM in our experiments. By default, we use BM25 (Robertson and Zaragoza, 2009) for sparse retrieval. In analysis, we additionally use Contriever (Izacard et al., 2022) for dense retrieval. During training, we use LoRA (Hu et al., 2022) for parameter-efficient fine-tuning, with Huggingface Accelerate for mixed-precision training. The entire training process takes approximately 3 hours on  $4 \times$  A100-SXM4-80GB GPUs.

<sup>4</sup>Details about benchmarks, retrieval, baselines, implementation, and task definition can be found in Appendix A.

Settings	2WikiMQA	HotpotQA	MuSiQue
<b>SGIR</b>	<b>74.47</b>	<b>63.09</b>	<b>37.15</b>
<i>Fine-grained Reward</i>			
(a) w/o. guided search	72.37	60.86	35.47
(b) w/o. reward	67.23	54.94	29.45
(c) w/o. self-critique	66.96	56.86	28.62
<i>Coarse-grained Reward</i>			
(d) w. coarse reward	73.40	59.85	34.32
(e) w/o. guided search	72.32	60.89	35.40

Table 2: Ablation experiments on rewards and search algorithm in fine-grained and coarse-grained scenarios. The experiments are conducted with Mistral-7B.

## 5 Experimental Results

We conduct extensive experiments on three multi-hop reasoning datasets, as shown in Table 1. We notice that although CoT is considered proficient in multi-step reasoning, but its performance is poor due to a lack of question-relevant knowledge within the LLM. When single-round retrieval is employed, an average performance improvement of 35.0% is observed. Iterative multi-round retrieval results in improved performance, achieving up to a 48.8% increase on 2WikiMQA. Additionally, methods based on question decomposition further improve reasoning performance. These results suggest that precise retrieval and reasoning planning are key factors affecting multi-hop reasoning.

Some methods equip the model with iterative reasoning capabilities through fine-tuning, introducing strategies such as active retrieval, question decomposition, and search during the iterative reasoning process. These methods can better address complex problems and achieve more accurate retrieval, further enhancing the performance.

As shown in Table 1, our method demonstrates remarkable superiority over baselines. It surpasses the previous state-of-the-art methods on all three datasets, achieving improvement of 2WikiMQA (+3.96), HotpotQA (+3.09), and MuSiQue (+4.85). Furthermore, our method excels particularly on the challenging MuSiQue dataset, which involves 3 to 4 hop hard questions. It achieves a significant improvement of 15.0% on the MuSiQue dataset compared to previous methods, with enhancements of 47.1% and 24.37% on 3/4-hop challenging questions, respectively. Additionally, we also conduct experiments on the Qwen2.5 and LLaMA2 models, achieving competitive performance that demonstrates the generalizability of our method.

We identify three key factors that contribute to

Stage	#Labeled	#Unlabeled	#Total	Quantile
Training Set	49000	-	-	-
Full SFT	38625	0	38625	100%
SP Start Point	15000	0	15000	95.3%
First Iteration	15000	21669	36669	97.1%
Second Iteration	15000	29659	44659	97.8%

Table 3: Comparison of the amount of data used for training and performance quantiles. #Labeled refers to data annotated by the teacher model with reasoning trajectories, and #Unlabeled refers to data synthesized by the smaller model itself. Quantile indicates the percentage of performance compared to full SFT.

our improvement: (a) Iterative multi-round reasoning and sub-question planning reduce the complexity of questions while improving retrieval accuracy. (b) Self-critique provides valuable feedback signals for intermediate steps during the reasoning process. (c) Reward-guided search selects superior reasoning trajectories among the exploration space.

## 6 Analysis

### 6.1 Ablation Study

We conduct ablation experiments to investigate the effects of different modules, as shown in Table 2. Our method employs fine-grained rewards by default. First, removing reward-guided search (a) and substituting it with greedy inference leads to considerable performance declines on all three datasets, highlighting the importance of search algorithms. Next, we remove the reward signal from greedy inference (b), which means the model’s inference process would not receive any guidance. This change results in a substantial performance decrease of 10%, demonstrating the essential role of reward signals in the inference process. Additionally, the exclusion of self-critique data from the training data produces a minor performance decline.

To further investigate the role of reward signals, particularly the impact of reward granularity on model performance, we replace the rewards with coarse-grained binary rewards, as shown in Table 2. In reward-guided greedy inference (e), fine-grained and coarse-grained rewards have similar effects. However, when a broader search is required (d), coarse-grained rewards exhibit certain drawbacks, resulting in a 2.7% performance decrease compared to fine-grained rewards. This indicates that fine-grained rewards are more effective in search-form reasoning, particularly when there are a substantial number of intermediate reasoning steps.

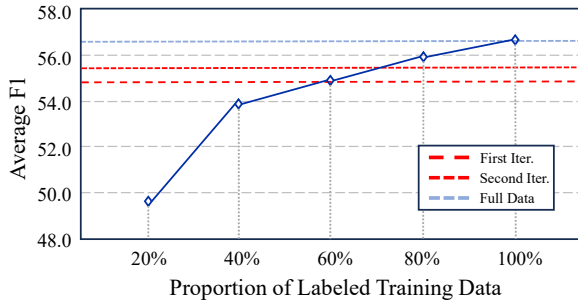


Figure 3: The impact of increasing training data and self-improvement on performance. As data scales, the performance consistently improves. Starting with 40% of the data, two self-improvement iterations boost performance to nearly 80%, with only 2% gap with full sft.

Strategy	2WikiMQA	HotpotQA	MuSiQue	Average
<b>Search Strategy</b>				
Guided Search	74.47	63.09	37.15	58.23
(a) one candidate	72.97	61.27	36.21	56.81
(b) no branching	68.87	52.88	28.34	50.03
Greedy Inference	72.37	60.86	35.47	56.23
<b>Selection Strategy</b>				
(c) cumulative reward	74.47	63.09	37.15	58.23
(d) outcome reward	71.79	60.96	35.23	55.99
(e) random select	71.17	60.10	35.17	55.81

Table 4: Analysis on various search and voting strategies on our method. Retaining more candidate reasoning trajectories and exploring different sub-questions help improve performance. Utilizing cumulative process rewards to select the final reasoning trajectory is effective.

## 6.2 Training Scaling and Self-Improvement

We investigate the performance under different amounts of training data and the effect of self-improvement data synthesis, as shown in Figure 3. Initially, as the amount of labeled data increases, the performance improves consistently. The most significant improvement occurs when 40% of the training data is utilized; beyond this threshold, the rate of performance enhancement decreases.

Starting with the model trained on 40% of the data, we perform iterative self-improvement. As shown in Figure 3 and Table 3, in the first iteration of self-improvement, the model achieves performance corresponding to the 60% training data quantile. In the second iteration, the model’s performance approaches the 80% quantile, with only a 2.2% gap from the performance achieved with the full labeled data fine-tuning. This suggests that our approach can effectively synthesize data through self-critique, yielding competitive performance in situations where labeled data is scarce.

Reward / Task	Retrieval	Reasoning	Overall	Average
Direct Reward	93.2	<b>95.8</b>	81.2	90.1
Post Explanation	95.4	95.4	80.6	90.5
Chain-of-thought	<b>97.4</b>	88.6	<b>87.4</b>	<b>91.1</b>
- Dedicated RM	95.8	93.6	78.8	89.4

Table 5: The performance of different reward forms on reward modeling tasks. Generative rewards are more effective than discriminative rewards, and mixed reward modeling has advantages over dedicated reward models in complex outcome reward modeling task.

Strategy	2WikiMQA	HotpotQA	MuSiQue	Average	Throughput
GenCritic	<b>72.37</b>	<b>60.86</b>	35.47	<b>56.23</b>	<b>1.92 QPS</b>
Separated	72.23	58.45	<b>36.46</b>	55.71	0.43 QPS

Table 6: Comparison of performance and efficiency between generative and separated critic model. The generative critic model brings slightly better performance, but has significant higher inference throughput.

## 6.3 Search and Trajectory Selection

We conduct experiments to investigate the factors influencing the guided reasoning process, including search strategies and reasoning selection, as shown in Table 4. As indicated in the preceding discussion, reward-based search is more effective than greedy inference. During the search process, we retain only one candidate trajectory in each iteration (a), which results in a performance decrease of 2.4%. This indicates that maintaining more candidate sequences preserves a greater number of potentially correct reasoning trajectories. Next, we set the question decomposition exploration width to 1 (b), which means only retrieval will be expanded during reasoning. This leads to a substantial drop of 13.57%, indicating the significance of exploring various sub-questions in the reasoning process.

At the end of the search, we need to select one trajectory as the final answer, and we investigate various answer selection strategies. Using cumulative process rewards for trajectory selection yields the best performance, followed by using outcome rewards. Both strategies outperform random selection. This suggests that in complex reasoning processes, intermediate rewards can serve as a criterion for assessing the reasoning quality, and that process rewards are more effective at distinguishing between different levels of reasoning quality.

## 6.4 Reward Modeling

The forms of reward (critic) models are diverse, such as generative and discriminative forms. We

Retriever	2WikiMQA	HotpotQA	MuSiQue	Average
<b>Reward-guided Greedy Inference</b>				
Sparse	72.37	60.86	35.47	56.23
Dense	75.44	59.53	38.48	57.81
Hybrid	<b>76.04</b>	<b>62.89</b>	<b>39.38</b>	<b>59.43</b>
<b>Reward-guided Search</b>				
Sparse	74.47	63.09	37.15	58.23
Dense	75.72	61.72	<b>40.36</b>	<b>59.26</b>
Hybrid	<b>76.05</b>	<b>63.47</b>	37.05	58.85

Table 7: The impact of retrieval systems on reasoning. Accurate retrieval leads to better reasoning performance. Our method naturally supports hybrid retrieval (sparse + dense) retrieval, which facilitates better reasoning.

examine the performance of various types of critic models on our reward modeling task, as shown in Table 5. Overall, generative reward modeling outperforms discriminative modeling. Moreover, among generative forms, chain-of-thought (CoT) is preferred over post-hoc explanations. Notably, CoT excels at modeling complex outcome rewards because it systematically decomposes the entire reasoning trajectory for step-by-step analysis.

Additionally, we investigate dedicated RM (each model assigned to a specific task) in comparison to mixed RM (one model responsible for all tasks). The dedicated RM exhibits performance variations of -1.6% / +5.6% for retrieval and reasoning tasks, respectively. However, there is a significant performance decrease of 9.8% in modeling outcome rewards. Since outcome rewards assess the entire reasoning process by integrating various reward modeling tasks, this suggests that multi-task reward modeling of sub-tasks enhances the model’s ability to evaluate complex outcome rewards.

Furthermore, we highlight that generative critic models offer greater benefits than separated generators and critics. As shown in Table 6, the GenCritic outperforms the separated models in two of the three datasets, achieving an average performance improvement of 0.9%. Without sacrificing reasoning performance, the GenCritic achieves a throughput that is 4.46 times higher, demonstrating significantly better reasoning efficiency<sup>5</sup>.

## 6.5 Retrieval System and Hybrid Retrieval

Retrieval plays a critical role in knowledge reasoning. To investigate the effect of the retrieval system, we employ various retrieval systems in the experiments, as illustrated in Table 7. Dense retrieval

<sup>5</sup>Please refer to Appendix B.1 for details of analysis.

Retriever	2WikiMQA	HotpotQA	MuSiQue	Average
BM25	72.37	60.86	35.47	56.23
Contriever	75.44	59.53	38.48	57.81
+ Hybrid	76.04	62.89	39.38	59.43
BGE-large	74.36	64.27	40.58	59.74
+ Hybrid	75.85	64.25	41.03	60.38
ColBERT-v2	75.37	65.44	42.24	61.02
+ Hybrid	75.37	65.69	43.05	61.37

Table 8: Experiments on different dense retrievers and their hybrid retrieval variants. Hybrid retrieval strategies improve performance across various dense retrievers.

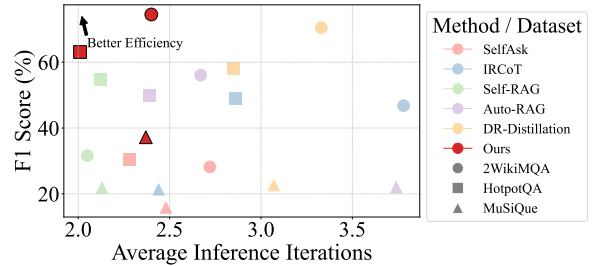


Figure 4: Relationship between performance and the average number of iterations to solve a multi-hop question. Upper-left represents efficient reasoning. Our method requires fewer iterations to solve one question with higher F1, demonstrating higher reasoning efficiency.

has better performance than sparse retrieval. When employing dense retrieval, a notable improvement of 2.8% is observed, suggesting that retrieval quality impacts reasoning performance. Furthermore, since our method treats each retrieved document independently, it naturally facilitates the integration of different retrieval systems. The hybrid retrieval system achieves a significant improvement of 5.7%, demonstrating the superiority of our method. However, in a guided search, excellent reasoning performance can be achieved using only sparse retrieval. This suggests that when retrieval quality is poor, the reward-guided search can help the model in exploring and filtering irrelevant retrieved documents, thereby improving reasoning performance.

In addition, we also apply hybrid retrieval on different dense retrievers (BGE-large<sup>6</sup> and ColBERT-v2<sup>7</sup>), as shown in Table 8. The hybrid retrieval strategy consistently achieves significant performance improvements across different dense retrievers, demonstrating its effectiveness.

<sup>6</sup><http://huggingface.co/BAAI/bge-large-en-v1.5>

<sup>7</sup><http://huggingface.co/colbert-ir/colbertv2.0>



## 6.6 Iterative Reasoning Efficiency

We assess reasoning efficiency by measuring the iterations needed to solve a multi-hop problem alongside its reasoning performance, as illustrated in Figure 4. Among the mentioned methods, Auto-RAG / DR-Distillation shows commendable reasoning capabilities, but it demands more iterations of reasoning, which incurs substantial costs. Conversely, our method not only surpasses their performance but also reduces the average number of reasoning iterations by 26.7%, demonstrating superior reasoning efficiency. Additionally, we also measure the average inference latency of single samples and token consumption in Appendix B.2.

## 7 Conclusion

In this work, we introduce Self-Critique Guided Iterative Reasoning for knowledge-intensive multi-hop question answering. Initially, we enable the model to conduct iterative reasoning through iterative question decomposition and to self-evaluate the intermediate reasoning steps, through a novel data synthesis and training pipeline. During inference, the model performs self-critique guided search: it first explores multiple reasoning paths via branching exploration and then utilizes self-critique feedback to refine the search process and identify the most promising reasoning trajectory. Extensive experiments and thorough analysis demonstrate the effectiveness of our proposed method.

### Limitations

Although we have conducted comprehensive experiments to demonstrate the effectiveness of our proposed methods, and provided detailed analysis on the components of our system, this paper still has some limitations.

For sub-question exploration, we employ a sampling-based approach rather than training a dedicated question decomposition module. Although this design improves the simplicity and efficiency of the overall system, it can limit the exploratory breadth of the model with respect to sub-questions. We plan to address this in future research, such as incorporating a query rewriting module.

Our method adopts a route of synthetic data followed by supervised fine-tuning (SFT), without incorporating reinforcement learning in the model’s sub-question reasoning process. In the future, we consider using strategies such as Monte

Carlo tree search to generate richer reasoning trajectories and introducing reinforcement learning to optimize these trajectories, thereby further enhancing the model’s reasoning performance.

### Acknowledgements

The research in this article is supported by the National Key Research and Development Project (2022YFF0903301), the National Science Foundation of China (U22B2059, 62276083). We also appreciate the support from China Mobile Group Heilongjiang Co., Ltd. @ on our research, the research is jointly completed by both parties. Ming Liu is the corresponding author.

### References

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. [Self-rag: Learning to retrieve, generate, and critique through self-reflection](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2022. [Improving language models by retrieving from trillions of tokens](#). In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 2206–2240. PMLR.
- Bradley C. A. Brown, Jordan Juravsky, Ryan Saul Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. 2024. [Large language monkeys: Scaling inference compute with repeated sampling](#). *CoRR*, abs/2407.21787.
- Shulin Cao, Jiajie Zhang, Jiaxin Shi, Xin Lv, Zijun Yao, Qi Tian, Lei Hou, and Juanzi Li. 2023. [Probabilistic tree-of-thought reasoning for answering knowledge-intensive complex questions](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 12541–12560. Association for Computational Linguistics.
- Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. 2024. [RQ-RAG: learning to refine queries for retrieval augmented generation](#). *CoRR*, abs/2404.00610.
- Zheng Chu, Jingchang Chen, Qianglong Chen, Haotian Wang, Kun Zhu, Xiyuan Du, Weijiang Yu, Ming

- Liu, and Bing Qin. 2024a. [Beamaggr: Beam aggregation reasoning over multi-source knowledge for multi-hop question answering](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 1229–1248. Association for Computational Linguistics.
- Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu, Bing Qin, and Ting Liu. 2024b. [Navigate through enigmatic labyrinth A survey of chain of thought reasoning: Advances, frontiers and future](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 1173–1203. Association for Computational Linguistics.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, Hao Zhang, Hanwei Xu, Hao Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang Yuan, Junjie Qiu, Junxiao Song, Kai Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Size Zheng, Tao Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaosha Chen, Xiaotao Nie, and Xiaowen Sun. 2024a. [Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model](#). *CoRR*, abs/2405.04434.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Size Zheng, Tao Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaosha Chen, Xiaotao Nie, and Xiaowen Sun. 2024b. [Deepseek-v3 technical report](#). *CoRR*, abs/2412.19437.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. 2023. [Retrieval-augmented generation for large language models: A survey](#). *CoRR*, abs/2312.10997.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. [Constructing A multi-hop QA dataset for comprehensive evaluation of reasoning steps](#). In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 6609–6625. International Committee on Computational Linguistics.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [Lora: Low-rank adaptation of large language models](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *CoRR*, abs/2311.05232.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. [Unsupervised dense information retrieval with contrastive learning](#). *Trans. Mach. Learn. Res.*, 2022.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth ee Lacroix, and William El Sayed. 2023a. [Mistral 7b](#). *CoRR*, abs/2310.06825.
- Jinhao Jiang, Jiayi Chen, Junyi Li, Ruiyang Ren, Shijie Wang, Wayne Xin Zhao, Yang Song, and Tao Zhang. 2024. [Rag-star: Enhancing deliberative reasoning with retrieval augmented verification and refinement](#). *CoRR*, abs/2412.12881.
- Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie

- Callan, and Graham Neubig. 2023b. [Active retrieval augmented generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 7969–7992. Association for Computational Linguistics.
- Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. [Retrieval-augmented generation for knowledge-intensive NLP tasks](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Xiang Li, Shizhu He, Fangyu Lei, JunYang JunYang, Tianhuang Su, Kang Liu, and Jun Zhao. 2024. [Teaching small language models to reason for knowledge-intensive multi-hop question answering](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 7804–7816. Association for Computational Linguistics.
- Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. [Query rewriting in retrieval-augmented large language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 5303–5315. Association for Computational Linguistics.
- OpenAI. 2023. [GPT-4 technical report](#). *CoRR*, abs/2303.08774.
- OpenAI. 2024. [Introducing openai o1](#).
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2023. [Measuring and narrowing the compositionality gap in language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 5687–5711. Association for Computational Linguistics.
- Zhenting Qi, Mingyuan Ma, Jiahang Xu, Li Lyna Zhang, Fan Yang, and Mao Yang. 2024. [Mutual reasoning makes smaller llms stronger problem-solvers](#). *CoRR*, abs/2408.06195.
- Stephen E. Robertson and Hugo Zaragoza. 2009. [The probabilistic relevance framework: BM25 and beyond](#). *Found. Trends Inf. Retr.*, 3(4):333–389.
- Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. 2023. [Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 9248–9274. Association for Computational Linguistics.
- Zhengliang Shi, Shuo Zhang, Weiwei Sun, Shen Gao, Pengjie Ren, Zhumin Chen, and Zhaochun Ren. 2024. [Generate-then-ground in retrieval-augmented generation for multi-hop question answering](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 7339–7353. Association for Computational Linguistics.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. [Scaling LLM test-time compute optimally can be more effective than scaling model parameters](#). *CoRR*, abs/2408.03314.
- Xin Su, Tiep Le, Steven Bethard, and Phillip Howard. 2024. [Semi-structured chain-of-thought: Integrating multiple sources of knowledge for improved language model reasoning](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, pages 8597–8613. Association for Computational Linguistics.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *CoRR*, abs/2307.09288.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. [Musique: Multi-hop questions via single-hop question composition](#). *Trans. Assoc. Comput. Linguistics*, 10:539–554.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. [Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 10014–10037. Association for Computational Linguistics.



Chaojie Wang, Yanchen Deng, Zhiyi Lv, Zeng Liang, Jujie He, Shuicheng Yan, and Bo An. 2024. [Q\\*: Improving multi-step reasoning for llms with deliberative planning](#). *CoRR*, abs/2406.14283.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.

Chunlei Xin, Yaojie Lu, Hongyu Lin, Shuheng Zhou, Huijia Zhu, Weiqiang Wang, Zhongyi Liu, Xianpei Han, and Le Sun. 2024. [Chain-of-rewrite: Aligning question and documents for open-domain question answering](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 1884–1896. Association for Computational Linguistics.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2024. [Qwen2.5 technical report](#). *CoRR*, abs/2412.15115.

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

Tian Yu, Shaolei Zhang, and Yang Feng. 2024. [Auto-rag: Autonomous retrieval-augmented generation for large language models](#). *CoRR*, abs/2411.19443.

Dan Zhang, Sining Zhou, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. 2024. [Rest-mcts\\*: LLM self-training via process reward guided tree search](#). In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.

## A Supplemental Materials on Setup

### A.1 Datasets

**HotpotQA** (Yang et al., 2018)<sup>8</sup> is a multi-hop reasoning dataset that includes two reasoning types: bridge and comparison. Bridge: sub-questions are linked by bridging entities. Comparison: compare two statements, including yes/no and partial order.

**2WikiMQA** (Ho et al., 2020)<sup>9</sup> is a multi-hop reasoning dataset that contains four reasoning types. The bridge (compositional) and comparison are the same as those in HotpotQA. Besides, it includes two additional reasoning types: bridge-comparison and inference. Bridge-comparison: compare two bridge sub-questions. Inference: sub-questions are connected through reasoning relationships rather than bridging entities. For example, the predicate grandfather is equal to a father’s father.

**MuSiQue** (Trivedi et al., 2022)<sup>10</sup> is a multi-hop reasoning dataset including 2-4 hops questions. All questions in MuSiQue are compositional, with the number of hops extended to a maximum of 4 hops. Besides, it includes multiple compositional questions, where a question contains more than one bridge entity. Due to the complexity of the questions, MuSiQue has the highest level of difficulty.

### A.2 Models

We use DeepSeek-V2.5 (DeepSeek-AI et al., 2024a)<sup>11</sup> as the advanced LLM to synthesis reasoning and critic rationales. In the experiments, we use Mistral-7B (Jiang et al., 2023a)<sup>12</sup> as the student backbone model by default. Besides, we also conduct experiments on LLaMA2-7B (Touvron et al., 2023)<sup>13</sup> and Qwen2.5-7B (Yang et al., 2024)<sup>14</sup> to demonstrate the generalization of our method.

### A.3 Training Corpus

We randomly sample subsets of the training data from 2WikiMQA, HotpotQA, and MuSiQue to construct the training corpus. Specially, we sample 20,000, 10,000, 19,000 instances from the respective datasets to construct the initial generator train-

<sup>8</sup><https://hotpotqa.github.io/>

<sup>9</sup><https://github.com/Alab-NII/2wikimultiHop>

<sup>10</sup><https://github.com/StonyBrookNLP/musique>

<sup>11</sup><https://huggingface.co/deepseek-ai/DeepSeek-V2.5>

<sup>12</sup><https://huggingface.co/mistral-community/Mistral-7B-v0.2>

<sup>13</sup><https://huggingface.co/meta-llama/Llama-2-7b>

<sup>14</sup><https://huggingface.co/Qwen/Qwen2.5-7B>



ing corpus  $X$ , while an additional 3,000, 2,000, and 1,000 non-overlapping instances are sampled as the critic model training corpus  $X_{critic}$ .

It is noteworthy that, to demonstrate the generality of our approach, we train a **unified** model on a mixture of all datasets, enabling it to perform reasoning across diverse multi-hop QA datasets instead of dataset-specific training.

#### A.4 Retrieval

Following prior work (Trivedi et al., 2023), we construct a retrieval corpus using supporting facts and irrelevant documents from the training instances, and conduct retrieval within this scope.

By default, we employ BM25 (Robertson and Zaragoza, 2009) provided by Elastic Search<sup>15</sup> for sparse retrieval in our experiments. Additionally, we incorporate dense retrieval and hybrid retrieval in our analysis. For dense retrieval, we use Contriever (Izacard et al., 2022)<sup>16</sup>. For hybrid retrieval, we combine the results from both sparse and dense retrieval by taking their union. Since our method independently processes each retrieved document and selects relevant ones based on self-evaluation, it inherently supports hybrid retrieval, thereby delivering better retrieval performance.

#### A.5 Baselines

**Chain-of-Thought** (Wei et al., 2022) generates step-by-step reasoning before the final answer. We use 20-shot demonstrations.

**RAG** (Lewis et al., 2020) retrieves  $k$  documents and conducts chain-of-thought reasoning conditional on the retrieval contents. We use BM25 to retrieve 5 documents, and 3-shot demonstrations.

**IRCoT** (Trivedi et al., 2023) interleaves retrieval and reasoning until the model arrives at the final answer. We set the number of documents per iteration to 4, with a maximum of 4 iterations.

**Self-Ask** (Press et al., 2023) iteratively breaks down complex questions and answers the resulting sub-questions one by one. We use 4-shot prompts.

**ProbTree** (Cao et al., 2023) decomposes the question into a tree and uses bottom-up aggregation reasoning based on log-prob to get the final answer.

<sup>15</sup><https://www.elastic.co/>

<sup>16</sup><https://huggingface.co/facebook/contriever-msmarco>

**BeamAggR** (Chu et al., 2024a) parses the complex question into a structured decomposition tree. It then performs bottom-up beam aggregation reasoning to select the promising reasoning trajectory.

**Self-RAG** (Asai et al., 2024) enables dynamic retrieval and generation through self-reflective tokens. We set the max retrieval documents per round to 5, and the max reasoning depth to 7.

**RAG-Star** (Jiang et al., 2024) uses Monte Carlo Tree Search to guide LLM-based tree-like reasoning and employs retrieval-augmented verification to consolidate internal and external knowledge.

**GenGround** (Shi et al., 2024) alternates between generating single-hop question-answer pairs using LLMs and grounding them in retrieved documents to correct potential knowledge errors.

**Auto-RAG** (Yu et al., 2024) leverages LLMs' reasoning and decision-making capabilities to plan retrievals, query rewriting, and gather external retrieval knowledge through multi-turn dialogues.

**DR-Distillation** (Li et al., 2024) distills two small models (Decomposer and Responder) to interactively solve knowledge-intensive multi-hop questions by breaking them into simpler sub-questions and conducting retrieval-augmented reasoning.

#### A.6 Task Definition

Given a multi-hop QA pair  $(Q, a)$  that requires multiple reasoning steps  $r_i$  to solve, each reasoning step  $r_i$  is based on a specific piece of supporting evidence  $d_i$ . We use the LLM to iteratively decompose the question  $q_i$ , employ an external retriever  $\text{Retr}$  to retrieve relevant documents based on the sub-question, and conduct retrieval-augmented reasoning conditional on the retrieved documents. Once the reasoning is finished, the answer in the last iteration is regarded as the final answer  $a^*$ .

$$Q \rightarrow r_1 \rightarrow r_2 \cdots \rightarrow a^* \quad (11)$$

#### A.7 Implementation Details

The training and inference are conducted on 4 and 1 NVIDIA Tesla A100-80G GPUs, respectively. The training process is initiated from the base model, and a chat template structured as '[INST] {Instruction} [/INST] {Output}' is manually added. During the training process, we exclusively train the Output while masking the retrieved document portions to prevent them from being included in the

Hyperparameters	Values
<i>Data Synthesis Hyperparameters</i>	
Temperature	1.15
N	3
# Raw Generator Corpus	49000
# Raw Critic Corpus	6000
<i>Training Hyperparameters</i>	
Epochs	2
Learning Rate	5e-5
Batch Size	128
Warmup Ratio	0.03
LR Decay	Cosine Annealing
LoRA Rank	128
LoRA Alpha	256
LoRA Targets	All Linear Layers
<i>Inference Hyperparameters</i>	
Temperature	1.15
Repetition Penalty	1.2
# Docs	5 or 6 (hybrid retrieval)
# Decomposition Branches	2
# Candidates	2
# Retrieval Reward	1.0 / 0.5 / -1.0
# Reasoning Reward	0.5 / 0.25 / -0.5

Table 9: Hyperparameters in experiments.

optimization. We use vLLM<sup>17</sup> in the inference. Hyperparameters are provided in Table 9.

<sup>17</sup><https://docs.vllm.ai/en/latest/>

## B Supplemental Materials on Analysis

Methods	2WikiMQA	HotpotQA	MuSiQue
Self-Ask	2.72	2.28	2.48
IRCoT	3.78	2.86	2.44
Self-RAG	2.05	2.12	2.13
Auto-RAG	2.67	2.39	3.74
DR-Distillation	3.33	2.85	3.07
SiGIR	<b>2.40</b>	<b>2.01</b>	<b>2.37</b>

Table 10: The average number of iterations of iterative reasoning methods to solve a multi-hop question. Our method requires fewer iterations to solve a multi-hop question, resulting in higher reasoning efficiency.

Question Type	Bdg.	Infer.	Cmp.	Bdg-Cmp.
Latency(s)	2.41	2.56	2.91	5.56
Prompt	2989.4	2775.0	1899.1	4653.9
Completion	171.5	202.4	209.3	375.2

Table 11: Per-sample inference latency and token consumption of greedy inference on 2WikiMQA dataset.

Question Type	Bridge	Comparison
Latency(s)	3.71	2.81
Prompt	2156.6	1786.6
Completion	182.7	220.3

Table 12: Per-sample inference latency and token consumption of greedy inference on HotpotQA dataset.

### B.1 Reasoning Efficiency between Separated Critic and Generative Critic

As mentioned in §6.4, GenCritic achieves a higher throughput compared to Separated Critic. We will describe how to calculate the throughput in this section. We use the same computing platform, equipped with one A100-SXM4-80GB GPU, a Xeon(R) 6348 CPU with 28 cores and 56 threads, and 480GB of memory. We sample 500 instances from each of the three datasets, totaling 1,500 instances. We conduct parallel inference using the API server provided by vLLM, employing reward-guided greedy inference during benchmarking. The throughput is calculated based on the time consumed for different methods.

### B.2 Iterative Reasoning Efficiency

We discuss the iterative reasoning efficiency in §6.6. Specifically, we consider a complete round of retrieval-generation as one reasoning iteration. It is

Question Type	2hop	3hop	4hop
Latency(s)	2.43	3.41	4.67
Prompt	2388.9	4162.6	4294.4
Completion	182.1	258.2	333.8

Table 13: Per-sample inference latency and token consumption of greedy inference on MuSiQue dataset.

noteworthy that multiple retrieval calls can occur within a single reasoning iteration. The specific number of reasoning iterations corresponding to each method can be found in Table 10. Additionally, we present the inference latency of individual samples with greedy inference, along with the average prompt and completion token costs, as shown in Table 11, Table 12, and Table 13.

Notation	Description
<b>(a) Modules</b>	
$LLM$	The advanced large language model to synthesis iterative reasoning rationales and process supervision.
$R$	The iterative reasoner model.
$C$	The separated critic model.
$R_{sc}$	The iterative reasoner model with self-critique capabilities.
Retr	External retriever (sparse, dense, hybrid)
<b>(b) Training Corpus</b>	
$X$	Raw training corpus for generator.
$X_{critic}$	Raw training corpus for critic, which is disjoint with $X$ .
<b>(c) Symbols</b>	
$q_i^{(r)}$	Multi-hop question after reduction.
$q_i$	Decomposed atomic sub-question.
$d_i$	The retrieved document.
$x$	Reasoning trajectory.
$s_i$	Sub-question reasoning.
$r_i^{retr}$	Reward for retrieval relevance.
$r_i^{reas}$	Reward for reasoning utility.
$r_i^c$	Cumulative process reward.
$r_i^o$	Outcome reward.

Table 14: The formal notations mentioned in Section 3.



## C Prompts

We provide the manually annotated prompts for data synthesis. For iterative reasoning rationales, we use 1-shot prompts for each question type (Figure 5). For critic on retrieval and reasoning, we use 3-shot prompts (Figure 6, 7). For critic on outcome reward, we use zero-shot prompts (Figure 8).

Additionally, we provide an example of iterative reasoning with self-critique, shown in Figure 15.

---

**Algorithm 1** Self-Critique Guided Iterative Reasoning

---

**Require:** Initial question  $q_0$ , retrieval model Retr, probabilistic language model  $p$ , candidate size  $k$ **Ensure:** Final answer  $a^*$ **Initialize:**

- 1:  $\mathbf{x}_0 \leftarrow [q_0]$  ▷ Initial reasoning prefix
- 2:  $r_0^c \leftarrow 0$  ▷ Initial c reward
- 3:  $\mathbf{x}_0 \leftarrow \{(\mathbf{x}_0, r_0^c)\}$  ▷ Set of candidate trajectories

**Main Loop:**

- 4: **for**  $t = 0$  to  $T - 1$  **do** ▷  $T$  is the maximum number of iterations
- 5:      $\mathbf{x}_{t+1} \leftarrow \emptyset$  ▷ Initialize set for next iteration

**Branch Exploration:**

- 6:     **for**  $(\mathbf{x}_t, r_t^c)$  in  $\mathbf{x}_t$  **do**
- 7:         **Decomposition:**
- 8:              $q_{t+1} \sim p(y \mid q_t^{(r)}, \mathbf{x}_t)$  ▷ Sample sub-question
- 9:         **Retrieval:R**
- 10:              $d_i \leftarrow \text{Retr}(q_{t+1})$  ▷ Retrieve relevant documents
- 11:              $r_{t+1}^{\text{retr}} \leftarrow p(y \mid q_{t+1}, d_i)$  ▷ Compute retrieval reward
- 12:         **Reasoning:**
- 13:              $s_{t+1} \sim p(y \mid q_{t+1}, d_i, \mathbf{x}_t)$  ▷ Sample reasoning step
- 14:              $r_{t+1}^{\text{reas}} \leftarrow p(y \mid q_{t+1}, d_i, s_{t+1})$  ▷ Compute reasoning reward
- 15:         **Update Reasoning Prefix:**
- 16:              $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t \oplus [q_{t+1}, d_i, r_{t+1}^{\text{retr}}, s_{t+1}, r_{t+1}^{\text{reas}}]$  ▷ Concatenate current status
- 17:         **Update Cumulative Reward:**
- 18:              $r_{t+1}^c \leftarrow r_t^c + r_{t+1}^{\text{retr}} + r_{t+1}^{\text{reas}}$
- 19:         **Question Reduction:**
- 20:              $q_t^{(r)} \leftarrow \text{ReduceQuestion}(\mathbf{x}_t)$  ▷ Reduce remaining question
- 21:              $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_{t+1} \cup \{(\mathbf{x}_{t+1}, r_{t+1}^c)\}$  ▷ Add new candidate
- 22:     **end for**

**Candidate Selection:**

- 23:      $\mathbf{x}_{t+1} \leftarrow \text{Top-k}(\mathbf{x}_{t+1}, r_{t+1}^c)$  ▷ Select top-k candidates
- 24: **end for**

**Final Answer Selection:**

- 25:  $a^* \leftarrow \arg \max_a (\sum_{i=1}^n (r_i^{\text{retr}} + r_i^{\text{reas}}))$  ▷ Select answer with highest c reward
  - 26: **return**  $a^*$
-

---

**Algorithm 2** Training SC-Reasoner

---

**Require:** Generator corpus  $X = \{x_i\}$ , where  $x_i = (q_i, a_i, \{d_1, \dots, d_n\})$

**Require:** Non-overlapping critic corpus  $X_{critic}$

**Require:** Large Language Model (LLM)

**Ensure:** Self-Critique Reasoner  $R_{sc}$

**Phase 1: Train Iterative Reasoner  $R$**

1:  $X_{synthetic} \leftarrow \emptyset$  ▷ Initialize synthetic reasoning corpus

2: **for** each instance  $x_i \in X$  **do**

3:   **Step 1.1: Synthesize Reasoning Rationales**

4:    $decomp_i, reasoning_i \leftarrow \text{LLM}(q_i, \{d_1, \dots, d_n\})$  ▷ Few-shot prompting for decomposition and reasoning

5:    $trajectory_i \leftarrow \text{Organize}(decomp_i, reasoning_i)$  ▷ Organize rationales

6:    $X_{synthetic} \leftarrow X_{synthetic} \cup \{trajectory_i\}$

7: **end for**

8: **Step 1.2: Train Iterative Reasoner**

9:  $R \leftarrow \text{TrainModel}(X_{synthetic})$  ▷ Train  $R$  on synthetic corpus

10:

**Phase 2: Train Critic Model  $C$**

11:  $X_{critic-data} \leftarrow \emptyset$  ▷ Initialize critic training data

12: **for** each instance  $x_i \in X_{critic}$  **do**

13:   **Step 2.1: Generate Reasoning Trajectories**

14:    $trajectory_i \leftarrow R(q_i, \{d_1, \dots, d_n\})$  ▷ Use  $R$  to generate reasoning trajectories

15:   **Step 2.2: Evaluate Rewards with LLM**

16:    $r_i^{retr}, r_i^{reas}, r_i^{overall} \leftarrow \text{LLM}(trajectory_i)$  ▷ Evaluate retrieval, reasoning, and overall quality

17:    $X_{critic-data} \leftarrow X_{critic-data} \cup \{(trajectory_i, r_i^{retr}, r_i^{reas}, r_i^{overall})\}$

18: **end for**

19: **Step 2.3: Train Critic Model**

20:  $C \leftarrow \text{TrainModel}(X_{critic-data})$  ▷ Train  $C$  on critic data

21:

**Phase 3: Train Self-Critique Reasoner  $R_{sc}$**

22:  $X_{sc-data} \leftarrow \emptyset$  ▷ Initialize self-critique training data

23: **for** each instance  $x_i \in X$  **do**

24:   **Step 3.1: Annotate Reasoning Trajectories with Critic**

25:    $trajectory_i \leftarrow R(q_i, \{d_1, \dots, d_n\})$  ▷ Generate reasoning trajectories with  $R$

26:    $r_i^{retr}, r_i^{reas}, r_i^{overall} \leftarrow C(trajectory_i)$  ▷ Annotate with critic  $C$

27:    $trajectory_i^{annotated} \leftarrow \text{AddSpecialTokens}(trajectory_i, r_i^{retr}, r_i^{reas}, r_i^{overall})$

28:    $X_{sc-data} \leftarrow X_{sc-data} \cup \{trajectory_i^{annotated}\}$

29: **end for**

30: **Step 3.2: Train Self-Critique Reasoner**

31:  $R_{sc} \leftarrow \text{TrainModel}(X_{sc-data})$  ▷ Train  $R_{sc}$  on annotated corpus

32:

**Self-Improvement Pipeline**

33: **while** not converged **do**

34:   **Step 4.1: Generate Synthetic Rationales**

35:    $X_{new} \leftarrow R_{sc}(\text{unlabeled questions})$  ▷ Generate rationales for unlabeled questions

36:   **Step 4.2: Evaluate Rationale Quality**

37:    $X_{high-quality} \leftarrow \{x \in X_{new} \mid C(x) > \text{threshold}\}$  ▷ Filter high-quality rationales

38:   **Step 4.3: Retrain  $R_{sc}$**

39:    $R_{sc} \leftarrow \text{TrainModel}(R_{sc}, X_{high-quality})$  ▷ Retrain with high-quality data

40: **end while**

41: **return**  $R_{sc}$

---

---

### Iterative Reasoning Examples of Self-Critique Reasoner

---

**Q** How many counties are there in the largest state in the region where the film The Family Stone was set? (Answer: sixteen)

---

[Non-Atomic Question]

<sub-question> In which region is the film "The Family Stone" set?

<paragraph>The Family Stone ...</paragraph>

Retrieval Quality: [Partially Relevant]

Reasoning: From Document, we know that "The Family Stone" is set in a small New England town. Therefore, the region is New England.

</sub-question> Reasoning Quality: [Fully supported]

[Remaining Question] How many counties are there in the largest state in New England?

[Non-Atomic Question]

<sub-question> Which state is the largest in New England?

<paragraph>New England ...</paragraph>

**A** Retrieval Quality: [Partially Relevant]

Reasoning: From Document, we know that Maine is the largest state in New England, constituting nearly one-half of the total area of New England.

</sub-question> Reasoning Quality: [Fully supported]

[Remaining Question] How many counties are there in Maine?

[Atomic Question]

<paragraph>List of counties in Maine ...</paragraph>

Retrieval Quality: [Relevant]

Reasoning: From Document, we know that there are sixteen counties in Maine.

Reasoning Quality: [Fully supported]

[Final Answer] @@**sixteen**@@

Overall Reasoning Quality: [4]

---

Table 15: A reasoning example of self-critique reasoner. The data instance is from MuSiQue dataset.



### Prompts for Iterative Reasoning Rationales

Please answer the following multi-hop questions with explanation based on the given documents. Wrap the final answer with @@ @@. Here is an example:

### Document #1: Cuidado con las imitaciones

Cuidado con las imitaciones is a 1948 Argentine film directed by Luis Bayón Herrera.

### Document #2: Luis Bayón Herrera

Luis Bayón Herrera (23 September 1889 – 30 March 1956) was a Spanish film director and screenwriter who worked in Argentine film of the 1940s and 1950s. He was one of the most important directors of the golden age of Argentine cinema. Herrera was born in Bilbao, Spain. He directed some 40 different Argentine feature films and in the last few years of his career in the early 1950s he worked on Cuban production with films such as "La Habana me voy" in 1951. He died in Buenos Aires in 1956, aged 66."

### Question: What nationality is the director of film Cuidado Con Las Imitaciones?

To answer this multi-hop question, we need to solve the following sub-questions:

### Decompose the multi-hop question into sub-questions

1. Who is the director of the film Cuidado Con Las Imitaciones?
2. What is the nationality of that director?

Next, let's solve the sub-questions one by one.

### Sub-question: Who is the director of the film Cuidado Con Las Imitaciones?

From Document #1, we know that "Cuidado Con Las Imitaciones" is a Spanish film, and its director is José María Forqué.

### Remaining Question: What is the nationality of José María Forqué?

From Document #2, we know that José María Forqué was a Spanish director, which means his nationality is Spanish.

### Final Answer

@@Spanish@@

{SUPPORTING DOCUMENTS}

### Question: {MULTI-HOP QUESTION}

Figure 5: The fewshot prompts for iterative reasoning rationales synthesis.

## Prompts for Retrieval Relevance Evaluation

You will be given a question and an evidence document. Your task is to evaluate whether the evidence is relevant to the question and provide a clear explanation for your decision. Use the following scale to rate the relevance of the evidence:

[Relevant]: The evidence directly addresses or supports the question.

[Partially Relevant]: The evidence is somewhat related to the question but does not fully address it. For example, it may mention the core entity or topic of the question.

[Irrelevant]: The evidence does not relate to the question in any meaningful way.

### Question: In what country was the first Pan-African conference held?

### Evidence: First Pan-African Conference

The First Pan-African Conference was held in London from 23 to 25 July 1900 (just prior to the Paris Exhibition of 1900 in order to allow tourists of African descent to attend both events). Organized primarily by the Trinidadian barrister Henry Sylvester Williams, it took place in Westminster Town Hall (now Caxton Hall) and was attended by 37 delegates and about 10 other participants and observers from Africa, the West Indies, the US and the UK, including Samuel Coleridge Taylor (the youngest delegate), John Alcindor, Dadabhai Naoroji, John Archer, Henry Francis Downing, and W.E.B. Du Bois, with Bishop Alexander Walters of the AME Zion Church taking the chair.

Rating: [Relevant]

Explanation: The evidence directly answers the question by stating that the first Pan-African conference was held in London, which is in the United Kingdom. The details provided, such as the dates and location of the conference, further confirm the relevance of the evidence to the question.

### Question: What is the administrative territorial entity for Malta Township?

### Evidence: Minsk Region

Minsk Region or Minsk Voblasć or Minsk Oblast (, "Minskaja voblasć" ; , "Minskaja oblastj") is one of the regions of Belarus. Its administrative center is Minsk, although it is a separate administrative territorial entity of Belarus. As of 2011, the region's population is 1,411,500.

Rating: [Irrelevant]

Explanation: The evidence provided discusses the Minsk Region, which is a region in Belarus. The question, however, asks about the administrative territorial entity for Malta Township. Since the evidence pertains to Belarus and not Malta, it is not relevant to the question.

### Question: Did Bob Rosa produce any work for Patrick Love?

### Evidence: Bob Rosa

Bob Rosa is an American record producer that has worked with artists such as Madonna, Janet Jackson, Mariah Carey and Espen Lind. One of his first recording projects was working on the track "Planet Rock" by Afrika Bambaataa & the Soulsonic Force in 1982.

Rating: [Partially Relevant]

Explanation: The evidence mentions Bob Rosa's work with various artists, including Madonna, Janet Jackson, Mariah Carey, and Espen Lind. However, it does not specifically mention Patrick Love or any work Bob Rosa may have done for him. Therefore, the evidence is partially relevant to the question.

### Question: {}

### Evidence: {}

Figure 6: The fewshot prompts for retrieval quality critic.

## Prompts for Reasoning Utility Evaluation

You will be provided with a question, an evidence document, along with a response. Your job is to determine whether the response is supported by the evidence, and provide explanation for your decision. Use the following scale to rate the response:

[Fully supported] Most of the information in the response is supported by the evidence.

[Partially supported] Some of the information in the response is supported by the evidence, but there are some parts that are speculative.

[Not supported] The response is not supported by the evidence / The response does not provide useful reasoning.

### Question: In which country is the Financial Conduct Authority located?

### Evidence: Martin Wheatley

Martin Wheatley is a British financier, formerly managing director of the Consumer and Markets Business Unit of the Financial Services Authority in the UK, and is the former CEO of the Financial Conduct.

### Response: From Document, we know that Martin Wheatley is a British financier who formerly managed the Consumer and Markets Business Unit of the Financial Services Authority in the UK. This indicates that the Financial Conduct Authority is located in the United Kingdom.

Rating: [Fully supported]

Explanation: The response is fully supported by the evidence. The evidence clearly states that Martin Wheatley is a British financier who formerly managed the Consumer and Markets Business Unit of the Financial Services Authority in the UK and was the former CEO of the Financial Conduct Authority. The response accurately extracts this information and concludes that the Financial Conduct Authority is located in the United Kingdom, which is directly supported by the evidence provided.

### Question: What is the birthplace of Jane Siberry?

### Evidence: Gavin Bradley

Gavin Bradley is an award-winning Canadian singer-songwriter and producer based in Toronto who has worked with artists like Nelly Furtado, Tori Amos and Jane Siberry. Fusing acoustic and electronic elements, his work is identifiable for its signature "warm" piano sound and live strings mixed with filtered synthesizers and other electronic manipulations. Besides production, Bradley is a solo recording artist. His debut album 'Deep Freeze' was released on UMI Records in 2006.

### Response: The response is partially supported because the evidence mentions Gavin Bradley's Toronto connection but does not directly state Jane Siberry's birthplace. The inference about her possible connection to Toronto is speculative.

Rating: [Partially supported]

Explanation: The response correctly identifies that the evidence mentions Gavin Bradley's connection to Toronto but does not directly state Jane Siberry's birthplace. The inference that Jane Siberry might be connected to Toronto based on her collaboration with Gavin Bradley is speculative, as the evidence does not provide explicit information about her birthplace. Therefore, the response is partially supported.

### Question: {}

### Evidence: {}

### Response: {}

Figure 7: The fewshot prompts for reasoning utility critic.

### Prompts for Outcome Reward Evaluation

You will be provided with a question, along with a reasoning trajectory.

Your job is to determine whether the reasoning trajectory is useful for answering the question, and provide explanation for your decision.

Use the following scale to rate the reasoning:

[5]: The reasoning process is clear, logically structured, and well-supported by the evidence.

[4]: The reasoning process is mostly clear, partially logically structured, and supported by evidence, but may contain minor logical flaws.

[3]: The reasoning is somewhat unclear with noticeable flaws and uses limited or weak evidence.

[2]: The reasoning is flawed, lacking supporting evidence, which results in an incorrect conclusion.

[1]: The reasoning fails to provide a definitive answer.

### Question: {}

### Reasoning: {}

Rating:

Explanation:

Figure 8: The prompts for reasoning trajectories outcome reward evaluation.