Search and Refine During Think: Facilitating Knowledge Refinement for Improved Retrieval-Augmented Reasoning

Yaorui Shi^{1*}, Sihang Li^{1*}, Chang Wu¹, Zhiyuan Liu², Junfeng Fang², Hengxing Cai^{3†}, An Zhang¹, Xiang Wang^{1†},

¹ University of Science and Technology of China
² National University of Singapore
³ DP Technology
{yaoruishi, sihang0520, xiangwang1223}@gmail.com, caihengxing@dp.tech
* Equal contribution. † Corresponding author.

Abstract

Large language models have demonstrated impressive reasoning capabilities but are inherently limited by their knowledge reservoir. Retrieval-augmented reasoning mitigates this limitation by allowing LLMs to query external resources, but existing methods often retrieve irrelevant or noisy information, hindering accurate reasoning. In this paper, we propose AutoRefine, a reinforcement learning post-training framework that adopts a new "search-and-refine-during-think" paradigm. AutoRefine introduces explicit knowledge refinement steps between successive search calls, enabling the model to iteratively filter, distill, and organize evidence before generating an answer. Furthermore, we incorporate tailored retrieval-specific rewards alongside answer correctness rewards using group relative policy optimization. Experiments on single-hop and multi-hop QA benchmarks demonstrate that AutoRefine significantly outperforms existing approaches, particularly in complex, multi-hop reasoning scenarios. Detailed analysis shows that AutoRefine issues frequent, higher-quality searches and synthesizes evidence effectively. Code is available at https://github.com/syr-cn/AutoRefine.

1 Introduction

Large language models (LLMs) have shown impressive abilities in language understanding, planning, and problem solving [1–3]. Recent advances demonstrate that reinforcement learning (RL) [4] further enhances LLMs' reasoning capabilities [5, 6], especially in complex tasks such as mathematics and coding [7, 8]. However, the knowledge encoded in LLMs is inherently constrained by their training corpora, limiting their reasoning performance on tasks requiring up-to-date information [9, 10].

A common strategy to address this limitation is retrieval-augmented generation (RAG), which equips LLMs with retrieval tools to access external knowledge bases during question answering [11–13]. Widely-adopted RAG pipelines typically rely on supervised fine-tuning (SFT) to train LLMs to issue search queries and generate responses based on retrieved documents [14–16]. While SFT can be effective for training large models for search, it sometimes necessitates the construction of high-quality search paths, which incurs additional effort and resource overheads [17]. To address this, recent studies draw inspiration from RL-based post-training [5] and explore RL for retrieval-augmented reasoning, achieving excellent results by only evaluating final answer correctness without the need for

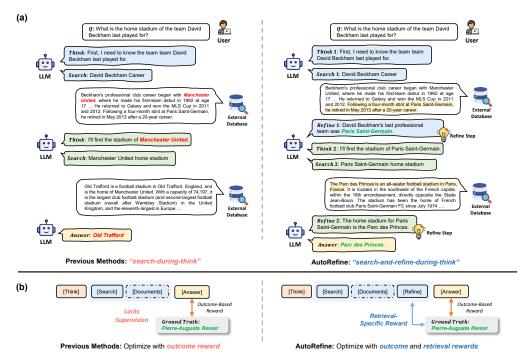


Figure 1: Comparison between previous retrieval-augmented reasoning methods and AutoRefine. (a) Previous "search-during-think" models can get distracted by irrelevant details between retrieval steps, leading to an incorrect answer. AutoRefine introduces a <Refine> step where the model explicitly refines crucial evidence, enabling the model to link information across multiple hops, plan its next query, and arrive at the correct final answer. (b) While previous methods rely only on outcome-based rewards, AutoRefine incorporates a retrieval-specific reward to directly supervise the <Refine> step.

pre-collected reasoning paths [18–23]. Scrutinizing existing studies on retrieval-augmented reasoning, we summarize a common "search-during-think" paradigm: given prompts with special search tokens (e.g., '<search> ... </search>'), the LLM is trained via RL to autonomously invoke retrieval tools, retrieve some documents from external knowledge bases, and generate answers within '<answer> ... </answer>' using the retrieved information.

Despite their promising results, we identify two core limitations inherent in the current retrievalaugmented reasoning paradigm:

- Lack refinement of retrieved documents. When facing out-of-scope questions, LLMs often require pieces of precise factual information (*e.g.*, names of historical figures, dates of events). However, the current "search-during-think" paradigm typically uses retrieval tools to return full documents based on input queries, many of which are noisy or only weakly relevant. As illustrated in Figure 1(a), previous methods reason directly over raw retrieved content, making it susceptible to getting distracted by irrelevant details. This is particularly problematic in multi-hop scenarios, where a distraction in an early step can derail the entire reasoning chain.
- Underexplored retrieval-specific rewards. While prior work on RL post-training highlights the importance of reward design [24–26], most retrieval-augmented reasoning methods rely solely on an outcome-based reward typically assessing the correctness of the final answer. As shown in Figure 1(b), this coarse supervision underexplores retrieval-specific rewards and offers little direct guidance for improving the retrieval process itself. As a result, it could be difficult for the LLM to learn how to retrieve more relevant or informative documents.

To address these limitations, we propose **AutoRefine**, a simple yet effective RL post-training framework that enhances the LLM's autonomous retrieval-augmented reasoning capability. At its core, AutoRefine adopts a "search-and-refine-during-think" paradigm, guided by a combination of answer and retrieval rewards. First, unlike prior "search-during-think" approaches that overlook refinement, we introduce an explicit knowledge refinement step into the reasoning loop using a 'search> ... </refine> ... </refine> template (§2.1). This template encourages the model

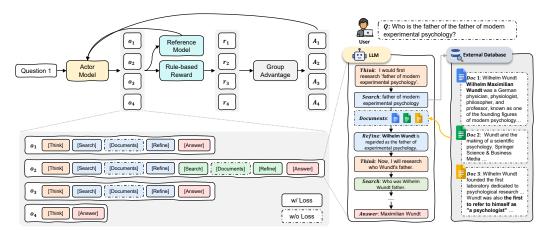


Figure 2: The training scheme of AutoRefine. (right) An actor model generates diverse reasoning trajectories for a given question, including think, search, refine, and answer. (left) These trajectories are optimized using the GRPO [7] algorithm described in Equation (4), where the loss on retrieved documents is masked out. We take G=4 in this example.

to explicitly distill crucial evidence from retrieved documents. By isolating key facts, the model can better link information across multiple retrieval hops and accurately plan subsequent queries before generating a final answer. Second, under this template, we apply Group Relative Policy Optimization (GRPO) [5] to train the model with both outcome-based and retrieval-specific rewards, rather than outcome-based rewards alone (§2.2). Specifically, during training, we first sample several trajectories from the model, each consisting of a sequence of <think>, <search>, <refine>, and <answer> steps. While the answer reward evaluates the final output, the retrieval reward is computed based on the quality of the content within the <refine> blocks, providing direct supervision for the refinement step. This joint reward design explicitly guides the model to extract, organize, and utilize fine-grained knowledge throughout reasoning.

To empirically assess AutoRefine, we conduct experiments on both single-hop [27–29] and multi-hop [30–33] question answering (QA) benchmarks. AutoRefine surpasses leading methods [5, 18, 21, 19] by 6.9% higher average accuracy, and shows especially high performances in multi-hop scenarios (cf. Table 1). It demonstrates a strong ability in identifying and addressing knowledge gaps via multi-turn, high-quality search queries. The knowledge refinement steps also effectively extract crucial information from noisy retrieved documents, directly contributing to improved answer quality. Additional experiments confirm the contribution of both the retrieval-specific reward and the refinement module, and AutoRefine holds robust performance across different retrieval depths.

2 Method

In this section, we introduce AutoRefine, a simple yet effective RL framework that enhances the LLM's autonomous retrieval-augmented reasoning capability. We first outline the overall task formulation and trajectory generation steps, highlighting the novel "search-and-refine-during-think" paradigm (§2.1). Next, we detail our reward modeling with both answer-based and retrieval-specific rewards to encourage fine-grained knowledge refinement (§2.2). Due to limited space, we refer to Appendix A for related work about reasoning in LLMs and retrieval augmented generation.

2.1 Trajectory Generation with Searching and Refinement

Task Formulation. Given a dataset $\mathcal{D} = \{(q, a)\}$ containing question–answer pairs and an external search engine \mathcal{E} , the task of retrieval-augmented reasoning requires the LLM to generate reasoning trajectories o by iteratively interacting with the knowledge source \mathcal{E} . Formally, for each question q, we generate a reasoning trajectory: $o = (\tau_1, \tau_2, \ldots, \tau_T)$, where the t-th intermediate reasoning step $\tau_t = (s_t, c_t)$ consists of an action $s_t \in \{\text{<think>}, \text{<search>}, \text{<documents>}, \text{<refine>}, \text{<answer>}\}$ and its associated content c_t . The model is expected to repeatedly retrieve and refine knowledge from \mathcal{E} until reaching a final answer o_{ans} that correctly addresses the question q.

You are a helpful assistant who is good at answering questions with multi-turn search engine calling. To answer questions, you must first reason through the available information using <think> and </think>. If you identify missing knowledge, you may issue a search request using <search> query </search> at any time. The retrieval system will provide you with the three most relevant documents enclosed in <documents> and </documents>. After each search, you need to summarize and refine the existing documents in <refine> and </refine>. You may send multiple search requests if needed. Once you have sufficient information, provide a concise final answer using <answer> and </answer>. <user> Question: {QUESTION} </user>

Figure 3: Prompt template for rollout generation.

Rollout Generation. The actor LLM π_{θ} generates trajectories by performing multiple rounds of interactions with the search engine \mathcal{E} . The trajectories contain multiple internal reasoning cycles, as illustrated in Figure 2. Each cycle consists of a sequence of structured operations: "<think>...</think>" for overall planning of consequent search actions, "<search>...</search>" for querying the external search engine, "<document>...</document>" for incorporating the retrieved documents, and "<refine>...</refine>" for distilling relevant information from the retrieved content. Following the reasoning phase, the model generates the final response within the "<answer>...</answer>" block based on the refined knowledge. Notably, the number of internal cycles is not manually pre-defined but autonomously determined by the actor LLM, adapting dynamically to the difficulty of the question. These tokens are defined and explained to the model via system instructions (cf. Figure 3).

Stopping Criteria. The generation terminates when an answer action is produced, *i.e.*, $s_T =$ <answer>. The content c_T of the terminal state τ_T is extracted as the trajectory's final answer o_{ans} .

2.2 Reward Modeling with Retrieval-Aware Signals

We use simple rule-based rewards to encourage free exploration during the RL process. The reward in AutoRefine consists of two complementary components: (1) the *Outcome-Based Reward*, which directly assesses the correctness of the answer generated by the model, and (2) the *Retrieval-Specific Reward*, which encourages the model to accurately identify and extract answer-relevant information from the retrieved documents.

Outcome-Based Reward. The outcome-based reward (aka). the answer reward) $\mathcal{R}_{Ans} \in [0,1]$ compares the model's final answer within the <answer></answer> block to the ground-truth answer, measuring its correctness. Formally, we treat the predicted and ground-truth answers as sets of words, and use F1-score between these two sets as the reward:

$$\mathcal{R}_{\text{Ans}} = \text{F1}(o_{\text{ans}}, a) = \frac{2|o_{\text{ans}} \cap a|}{|o_{\text{ans}}| + |a|},\tag{1}$$

where o_{ans} is the predicted answer (e.g., "Pierre-Auguste Renoir" in yellow box of Figure 1(a), and a is the ground truth answer from the (q, a) pair (e.g.), the gray box in Figure 1(b)).

Retrieval-Specific Reward. We further introduce an additional reward, the retrieval reward $\mathcal{R}_{Ret} \in \{0,1\}$, to explicitly encourage the extraction and utilization of relevant information from noisy retrieved documents. The retrieval reward is measured based on the quality of refined documents within the <refine></refine> blocks. Specifically, we collect all knowledge refinement steps (*i.e.*, content within the <refine>...</refine> blocks) across the trajectory and concatenate them into a single text sequence:

$$\mathcal{R}_{\text{Ret}} = \mathbb{I}(a \cap o_{\text{refine}} = a),$$
 (2)

where $\mathbb{I}(\cdot)$ is the indicator function, $o_{\text{refine}} = \bigcup \{c_t \mid (s_t, c_t) \in o \land s_t = \text{refine}\}\$ is the concatenation of all the knowledge refinement steps (e.g., "The documents concludes...'The Umbrellas'." in blue box of Figure 1). This reward is activated when all components of the ground-truth answer are present in the refined knowledge, reinforcing faithful and targeted information extraction.

Integrating Outcome and Retrieval Rewards. The overall reward function in AutoRefine is designed to encourage both accurate final answers and meaningful intermediate knowledge extraction. Specifically, the model receives a full reward of 1 if it generates the correct answer. If the final answer

is incorrect but some relevant information has been extracted during the refinement step, a partial reward of 0.1 is assigned. No reward is granted if neither correct answers nor relevant information are produced. The overall reward $\mathcal{R}_{\text{Overall}}$ can be formally written as:

$$\mathcal{R}_{\text{Overall}} = \begin{cases}
\mathcal{R}_{\text{Ans}}, & \text{if } \mathcal{R}_{\text{Ans}} > 0 \\
0.1, & \text{if } \mathcal{R}_{\text{Ans}} = 0 \text{ and } \mathcal{R}_{\text{Ret}} > 0 \\
0. & \text{if } \mathcal{R}_{\text{Ans}} = \mathcal{R}_{\text{Ret}} = 0
\end{cases}$$
(3)

Training Objective. We apply Group Relative Policy Optimization (GRPO) [7] as the policy optimization algorithm for RL. The overview of the GRPO training scheme is shown in the top-left corner of Figure 2. Formally, given an actor model π_{θ} and a reference model π_{ref} , a group of G rollouts $\{o_i\}_{i=1}^G$ is sampled as described in §2.1. We optimize the actor model π_{θ} by maximizing:

$$\underset{\theta}{\operatorname{argmax}} J_{\text{GRPO}}(\theta) = \mathbb{E}_{(q,a) \sim \mathcal{D}, \{o_{i}\}_{i=1}^{G} \sim \pi_{\theta_{\text{old}}}(\cdot | q)} \left[\frac{1}{G} \sum_{i=1}^{G} \frac{1}{|o_{i}|} \sum_{t=1}^{|o_{i}|} \min \left(\frac{\pi_{\theta}(o_{i,t} \mid q, o_{i, < t})}{\pi_{\theta_{\text{old}}}(o_{i,t} \mid q, o_{i, < t})} \hat{A}_{i,t}, \right. \\
\left. \text{clip} \left(\frac{\pi_{\theta}(o_{i,t} \mid q, o_{i, < t})}{\pi_{\theta_{\text{old}}}(o_{i,t} \mid q, o_{i, < t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right) - \beta \mathbb{D}_{\text{KL}} \left[\pi_{\theta} \parallel \pi_{\text{ref}} \right] \right] \tag{4}$$

where $\hat{A}_{i,t} = [r_{i,t} - \text{mean}(r_t)]/\text{std}(r_t)$ is the normalized token-level advantage for the i-th rollout in the group, G is the group size, ϵ is the clipping ratio, and β is the coefficient for the estimated KL divergence. As shown in the bottom-left part of Figure 2, we mask out the retrieved documents during the loss computation.

3 Experiments

In this section, we aim to answer the following Research Questions (RQs):

RQ1: How effectively does AutoRefine's "search-and-refine-during-think" paradigm enhance performance in retrieval-augmented question answering?

RQ2: Can AutoRefine effectively resolve information gaps through retrieval, especially when facing complex multi-hop problems?

RQ3: Can knowledge refinement steps distill critical information from retrieved documents?

RQ4: Can AutoRefine achieve robust performance improvements under different retriever settings?

3.1 Experiment Setup

Datasets. We evaluate performance using seven diverse QA benchmarks, including three single-hop QA datasets: Natural Questions (NQ) [27], TriviaQA [29], PopQA [28], and four datasets that require multi-hop searching: HotpotQA [30], 2WikiMultihopQA (2Wiki) [31], Musique [32], Bamboogle [33] for evaluation. Exact match accuracy serves as the evaluation metric for all downstream datasets. Following the setting of prior works [19], we train AutoRefine using a combined training set from NQ and HotpotQA.

Baselines. In our experiments, we compare AutoRefine against three kinds of methods: (1) generation without retrieval (w/o Retrieval), including direct generation with LLM, supervised fine-tuning (SFT), and R1-like training (R1) [5] without Retrieval; (2) methods with single-hop retrieval (w/ Single-Hop Retrieval), including direct retrieval with the input question (Naive RAG); (3) training with retrieval (w/ Multi-Hop Retrieval) including: agentic search method Search-o1 [18], IRCoT [34], retrieval-augmented reasoning model Search-R1 [19] and ReSearch [21].

Implementation Details. To simulate a real-world search scenario, we remove original context documents from the QA datasets [30–32] and instead use the December 2018 Wikipedia dump [35] as the external knowledge source, with E5-base-v2 [36] as the retrieval engine. By default, the search engine retrieves the top three most relevant documents on each query. For RL-based baselines, we run experiments using both Qwen2.5-3B-Base and -Instruct models. For SFT and direct generation

Table 1: (**RQ1**) Accuracy comparison of AutoRefine versus baseline methods with Qwen2.5-3B [3] across various QA benchmarks. **Bold** denotes best results, and <u>underline</u> denotes second best results.

		Single-Hop (Multi-Hop QA					
Methods	NQ	TriviaQA	PopQA	HotpotQA	2Wiki	Musique	Bamboogle	Avg.
w/o Retreival								
Direct Generation	0.106	0.288	0.108	0.149	0.244	0.020	0.024	0.134
SFT	0.249	0.292	0.104	0.186	0.248	0.044	0.112	0.176
R1-Instruct [5]	0.210	0.449	0.171	0.208	0.275	0.060	0.192	0.224
R1-Base [5]	0.226	0.455	0.173	0.201	0.268	0.055	0.224	0.229
w/ Single-Hop Retrieval								
Naive RAG [37]	0.348	0.544	0.387	0.255	0.226	0.047	0.080	0.270
w/ Multi-Hop Retrieval								
Search-o1 [18]	0.238	0.472	0.262	0.221	0.218	0.054	0.320	0.255
IRCoT [34]	0.111	0.312	0.200	0.164	0.171	0.067	0.240	0.181
ReSearch-Instruct [21]	0.365	0.571	0.395	0.351	0.272	0.095	0.266	0.331
ReSearch-Base [21]	0.427	0.597	0.430	0.305	0.272	0.074	0.128	0.319
Search-R1-Instruct [19]	0.397	0.565	0.391	0.331	0.310	0.124	0.232	0.336
Search-R1-Base [19]	0.421	0.583	0.413	0.297	0.274	0.066	0.128	0.312
AutoRefine-Instruct	0.436	0.597	0.447	0.404	0.380	0.169	0.336	0.396
AutoRefine-Base	0.467	0.620	0.450	0.405	0.393	<u>0.157</u>	0.344	0.405

baselines, we use the instruct variant to better align with instruction-following tasks. Most baseline results are taken from Search-R1 [19], which has experimental settings consistent with ours. We reproduce ReSearch using the authors' publicly available code. Additional implementation details can be found in Appendix B, and further experimental results are provided in Appendix C.

3.2 Overall Performance (RQ1)

Table 1 presents the overall performance comparison between AutoRefine and the baseline methods. The Avg. column stands for the average accuracy. As shown in the results, AutoRefine significantly outperforms baseline models across the seven benchmarks. It achieves a 0.069 accuracy gain on the base variant and a 0.060 improvement on the instruct one compared to the strongest baseline.

Besides the overall performance, we observe that the performance gains achieved by AutoRefine are more obvious on the multi-hop QA benchmarks. For example, AutoRefine improves the performance on 2Wiki by 0.083 and Musique by 0.045, which implies 21% and 26.7% relative increase, respectively. According to further analysis in §3.3.1 and §3.3.2, we attribute AutoRefine's extraordinary performance on multi-hop benchmarks to its ability to perform high-quality searching and efficient utilization of retrieved documents.

Obs 1: AutoRefine significantly improves QA accuracies, especially on multi-hop benchmarks.

3.3 Analytical Results

3.3.1 Search Behaviors (RQ2)

A crucial capability of retrieval-augmented reasoning models is identifying and addressing knowledge gaps via retrieval. To evaluate this, we analyze the **search frequency** and **search quality** of AutoRefine in four scenarios: training samples, all seven downstream benchmarks, single-hop QA benchmarks, and multi-hop QA benchmarks. These analyses are depicted in Figure 4.

Search Frequency reflects the model's capability to recognize knowledge gaps and perform searches accordingly. This search behavior is measured by the average number of search calls a model makes per rollout. We analyze the search frequency of both -Base and -Instruct variants of AutoRefine in Figure 4(a). As the figure shows, both variants have evolved multi-turn searching abilities after enough training steps. The average number of search calls converges to around 1.5 for AutoRefine-Instruct, and higher than 2 for AutoRefine-Base.

Despite the high overall search frequency on seven benchmarks, AutoRefine demonstrates distinct search behaviors for single-hop and multi-hop questions. On the three single-hop benchmarks, both

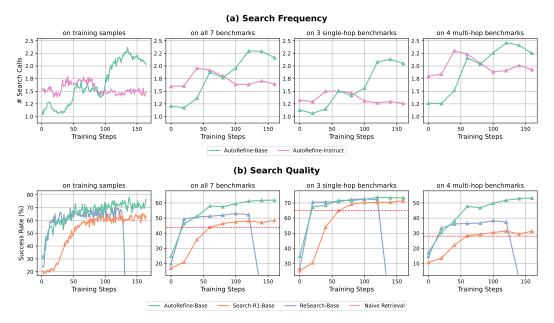


Figure 4: (**RQ2**) Visualization of the search behaviors. (a) AutoRefine's average number of search calls per rollout. For both variants, AutoRefine learns to adaptively issue more search queries for multi-hop questions and fewer for single-hop ones. (b) Comparison of search success rates between retrieval-augmented reasoning methods. While all methods draft more efficient search queries than naive retrieval, AutoRefine achieves more significant performance gains.

variants begin with fewer than 1.3 searches per rollout, and gradually adjust to 2.0 and 1.2. In contrast, the models begin with much higher search frequencies when facing multi-hop questions, which rapidly go up to $2.0 \sim 2.5$. This phenomenon exhibits AutoRefine's ability to dynamically adjust the number of search calls according to the complexity of downstream tasks, with more frequent searching on multi-hop questions and less on single-hop ones.

Obs 2.1: AutoRefine learns to perform multi-turn searching and can adaptively issue search queries depending on task complexity.

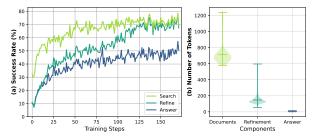
Search Quality evaluates whether the model generates effective search queries that can return informative documents. Knowledge-intensive questions often demand precise factual information to ask, e.g., names of historical figures or dates of events. In such cases, search calls can only be considered successful if the retrieved documents directly contain the answer. Hence, we estimate the search quality by counting the proportion of successful searches where retrieved documents contain the ground truth answer a. We also include naive retrieval, which directly uses the input question to conduct one-turn searches, as a reference.

We compare the search quality of AutoRefine against Search-R1 and ReSearch trained from Qwen2.5-3B-Base in Figure 4(b). All reasoning-based methods learns to draft efficient queries that have higher searching quality than naive retrieval. In single-hop scenarios, all three methods converge to a high success rate of about 70% after 100 steps. On multi-hop benchmarks, the search qualities of baseline methods also successfully converge to $30\% \sim 40\%$, much higher than that of baseline methods. In contrast, the search quality of AutoRefine continuously goes up to higher than 50%, which surpasses baseline methods by a large margin of $10\% \sim 15\%$.

Obs 2.2: AutoRefine drafts efficient queries that retrieve documents relevant to the answer.

3.3.2 Effectiveness of Knowledge Refinement (RQ3)

To investigate the effectiveness of knowledge refinement, we specifically analyze whether the refinement steps successfully distill critical information from retrieved documents. We start by comparing the success rates of different actions: <search>, <refine>, and <answer>. Here we use cover exact



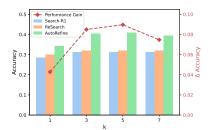


Figure 5: **(RQ3)** Comparison of search, refine, and answer actions over (a) recall and (b) average token counts. Knowledge refinement keeps crucial information from retrieved documents while reducing context length.

Figure 6: **(RQ4)** Comparison of downstream accuracies under different retrieval depths. AutoRefine exhibits robust gains for $1 \le k \le 7$.

match to measure the success rate for all three actions, which is defined as the proportion of actions that return documents/refinements/answers containing the ground truth answer.

The results are shown in Figure 5. After enough training steps, the success rate of <refine> actions tends to align with that of the <search> action (Figure 5(a)). This suggests the model gradually learns to keep crucial evidence as long as the search returns correct documents. Figure 5(b) also provides the length of each component. The token count of refinement steps is about $100 \sim 200$ tokens, which is about 4 times fewer than the documents ($\geqslant 600$ tokens). Comparing Figure 5(a) and (b), we find that the knowledge refinement steps of AutoRefine greatly reduce the context length, while successfully preserving the information that is relevant to the answer.

Obs 3: Knowledge refinement steps efficiently distill critical evidence from retrieved documents while filtering out irrelevant content.

3.3.3 Impact of Retrieval Depths (RQ4)

Different retriever settings may also influence retrieval-augmented generation models, and one important aspect is the retrieval depth. While more documents per search could potentially provide richer external knowledge, it also includes more noise in documents. To explore the models' robustness across different retrieval depths, we vary the number of documents (k) returned by the retrieval engine at evaluation time, from 1 to 7, while training fixed to k=3.

The comparison of inference accuracy is shown in Figure 6. All three methods achieve robust performance across different k levels. Compared to the baseline methods, AutoRefine steadily boosts the average accuracy by $0.04 \sim 0.1$, demonstrating its strong document denoising ability. The accuracy increments caused by AutoRefine are particularly obvious when $k \geqslant 3$, which is likely caused by its strong ability to discover useful information under increasingly noisy conditions. Peak performance gain is 0.09 observed at k=5, where a balanced trade-off between information richness and noise is reached.

Obs 4: AutoRefine exhibits consistent improvements across varying retrieval depths.

3.4 Ablation Studies

3.4.1 Ablation on Key Components

We conduct ablation studies over the key components in AutoRefine. Specifically, we consider three configurations: (1) the full AutoRefine model, (2) AutoRefine without the retrieval-specific reward \mathcal{R}_{Ret} (w/o Retrieval Reward), and (3) AutoRefine without both \mathcal{R}_{Ret} and the knowledge refinement step (w/o Retrieval Reward & Refinement). See Appendix C.3 for analysis of retrieval reward design.

Impact on Answer Accuracy. Table 2 presents the answer accuracy on downstream benchmarks for each configuration. The results demonstrate that both the retrieval-specific reward and the knowledge refinement step are essential for achieving strong performance. The full AutoRefine model consistently achieves the highest average accuracy across both the base and instruct variants.

Table 2: Ablation study over key components in AutoRefine.

		Single-Hop QA			Multi-Hop QA				
Model Variants	NQ	TriviaQA	PopQA	HotpotQA	2wiki	Musique	Bamboogle	Avg.	
AutoRefine-Base	0.467	0.620	0.450	0.405	0.393	0.157	0.344	0.405	
w/o Retrieval Reward	0.423	0.583	0.424	0.368	0.351	0.139	0.344	0.376	
w/o Retrieval Reward & Refinement	0.422	0.585	0.419	0.294	0.257	0.062	0.144	0.312	
AutoRefine-Instruct	0.436	0.597	0.447	0.404	0.380	0.169	0.336	0.396	
w/o Retrieval Reward	0.418	0.587	0.429	0.355	0.335	0.124	0.272	0.360	
w/o Retrieval Reward & Refinement	0.406	0.580	0.412	0.319	0.312	0.091	0.210	0.333	



Figure 7: Effectiveness of key components over the search behaviors and the refinement quality.

Impact on Search and Refinement Abilities. We further analyze how each component affects AutoRefine's search and refinement capabilities. Figure 7 shows a comparative analysis using Qwen2.5-3B-Base. As illustrated in Figure 7(a), the retrieval-specific reward effectively promotes multi-turn search behavior. Additionally, it significantly boosts knowledge refinement quality, yielding approximately a 20% improvement in refinement success rate (Figure 7(c)). The inclusion of the knowledge refinement step also enhances both the frequency and quality of retrieval, as shown in Figures 7(a) and (b).

3.4.2 Ablation on Model Sizes and Evaluation Metrics

To provide a comprehensive evaluation of AutoRefine, we conduct ablation studies on two aspects: (1) the model sizes, including Qwen2.5-3B and Qwen2.5-7B; (2) the evaluation metrics, including exact match (EM), F1 score, and cover exact match (CEM). The results are shown in Table 3.

Comparing the performance of AutoRefine on Qwen2.5-3B and Qwen2.5-7B, we observe that the larger model size generally leads to better performance. AutoRefine-7B achieves approximately 0.05 performance gains on all metrics, which is slightly lower compared to those on the 3B variant. Besides, AutoRefine maintains superior performance on all three metrics compared to baselines.

3.4.3 Ablation on Knowledge Refinement Module

To demonstrate the necessity of the RL-driven refinement steps, we compare AutoRefine with several baselines that use external summarization models as refiners. These baselines augment Search-R1 with refiners based on BART [38] and Qwen2.5-3B-Instruct [3]. For the Qwen model, we test two prompting strategies: one that only asks for summarization, and another that asks for both summarization and a plan for the next search step. The results are presented in Table 4.

The experiment indicates that simply adding an external summarizer to Search-R1 improves performance on some single-hop QA benchmarks (*e.g.*, PopQA) but can be detrimental in multi-hop settings. In contrast, AutoRefine maintains superior performance on hard multi-hop benchmarks. Through RL, AutoRefine learns not only to summarize but also to introspect, identify missing information, and plan its next actions We find out its performance gain on multi-hop benchmarks derives from the ability to not just summarize, but also to introspect, recognize missing information, and plan subsequent search steps, as illustrated in the case studies (§C.5).

Table 3: Ablation study over model sizes and evaluation metrics.

			General QA	1					
Model	Metric	NQ	TriviaQA	PopQA	HotpotQA	2wiki	Musique	Bamboogle	Avg.
Qwen2.5-7B	-Base								
	EM	0.469	0.627	0.449	0.410	0.272	0.173	0.456	0.408
Search-R1	F1	0.552	0.700	0.487	0.517	0.327	0.236	0.560	0.483
	CEM	0.509	0.680	0.467	0.445	0.309	0.197	0.496	0.443
	EM	0.484	0.659	0.487	0.451	0.405	0.187	0.512	0.455
AutoRefine	F1	0.574	0.729	0.525	0.573	0.467	0.283	0.604	0.536
	CEM	0.523	0.707	0.500	0.487	0.441	0.217	0.528	0.486
Qwen2.5-3B	-Base								
	EM	0.421	0.583	0.413	0.297	0.274	0.066	0.128	0.312
Search-R1	F1	0.476	0.650	0.429	0.380	0.322	0.123	0.184	0.366
	CEM	0.462	0.642	0.442	0.325	0.288	0.082	0.128	0.338
	EM	0.467	0.620	0.450	0.405	0.393	0.157	0.344	0.405
AutoRefine	F1	0.534	0.689	0.479	0.503	0.453	0.233	0.449	0.477
	CEM	0.502	0.674	0.468	0.440	0.428	0.175	0.384	0.439

Table 4: Performance comparison against Search-R1 with external refiners.

	General QA							
Model	NQ	TriviaQA	PopQA	HotpotQA	2wiki	Musique	Bamboogle	Avg.
AutoRefine	0.467	0.620	0.450	0.405	0.393	0.157	0.344	0.405
Search-R1	0.421	0.583	0.413	0.297	0.274	0.066	0.128	0.312
Search-R1 + Refiner (BART [38])	0.395	0.619	0.450	0.337	0.239	0.065	0.115	0.317
Search-R1 + Refiner (Qwen, Summary)	0.399	0.600	0.445	0.331	0.264	0.073	0.180	0.328
Search-R1 + Refiner (Qwen, Summary & Plan)	0.378	0.562	0.431	0.299	0.231	0.059	0.149	0.301

4 Limitations

Despite the promising performance of AutoRefine in retrieval-augmented reasoning tasks, several limitations remain for further investigation.

- Evaluation Metrics. This work evaluates model performance solely on exact match accuracy or F1 score, which may overlook semantically correct responses with minor textual variations. This limits the evaluation of long-form or open-ended responses.
- Static Retrieval Corpus. The retrieval component uses a fixed Wikipedia snapshot, lacking current or time-sensitive information. This setting limits the system's applicability to real-world use cases where users expect information from live search engines.

5 Conclusion and Future Work

This work proposes AutoRefine, an RL post-training framework designed to improve the retrieval-augmented reasoning capabilities of LLMs. AutoRefine adopts a novel "search-and-refine-during-think" paradigm that explicitly encourages the model to identify and distill relevant information from noisy retrieved content. By jointly optimizing for both outcome-level and retrieval-specific rewards, AutoRefine effectively guides LLMs to extract, assess, and integrate external knowledge. Comprehensive evaluations show that AutoRefine consistently surpasses existing methods, achieving up to a 6.9% average improved accuracy on seven QA benchmarks. These results underscore its potential to enhance the accuracy and reliability of retrieval-augmented LLMs.

Future work will focus on addressing the limitations identified above. First, we aim to adopt more flexible and semantically aware evaluation metrics — such as LLM-as-a-Judge evaluation — to more effectively measure answer quality in complex question answering tasks. Second, we intend to adapt AutoRefine to dynamic retrieval settings, including live web search and continuously evolving document corpora. By addressing these directions, we aim to further improve the sflexibility and time-sensitivity of AutoRefine, thus broadening its practicality in more realistic applications. This extension would enable the system to operate in more realistic, time-sensitive applications and broaden its practical utility.

6 Acknowledgement

This research is supported by the National Natural Science Foundation of China (62572449). We also appreciate the support of Siyuan Wang¹ from Shanghai Jiao Tong University.

https://scholar.google.com/citations?user=HqGFfDkAAAAJ&hl=en

References

- [1] OpenAI. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- [2] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [3] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [4] Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. *Journal of artificial intelligence research*, 4:237–285, 1996.
- [5] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, 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, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. CoRR, abs/2501.12948, 2025. doi: 10.48550/ARXIV.2501.12948.
- [6] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, and Ilge Akkaya. Openai o1 system card. CoRR, abs/2412.16720, 2024. doi: 10.48550/ARXIV.2412.16720.
- [7] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [8] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.
- [9] Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. Active retrieval augmented generation. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *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, 2023. doi: 10.18653/V1/2023. EMNLP-MAIN.495.
- [10] Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. Enhancing retrieval-augmented large language models with iterative retrieval-generation synergy. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 9248–9274. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.FINDINGS-EMNLP. 620.
- [11] Akari Asai, Sewon Min, Zexuan Zhong, and Danqi Chen. Retrieval-based language models and applications. In Yun-Nung Vivian Chen, Margot Mieskes, and Siva Reddy, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts, ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 41–46. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.ACL-TUTORIALS.6.
- [12] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *CoRR*, abs/2312.10997, 2023. doi: 10.48550/ARXIV.2312.10997.
- [13] Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. In-context retrieval-augmented language models. *Trans. Assoc. Comput. Linguistics*, 11:1316–1331, 2023. doi: 10.1162/TACL_A_00605.
- [14] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. Open-Review.net, 2024.
- [15] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. REPLUG: retrieval-augmented black-box language models. In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard, editors, 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 8371–8384. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.NAACL-LONG.463.
- [16] Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. Corrective retrieval augmented generation. *arXiv preprint arXiv:2401.15884*, 2024.
- [17] Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V. Le, Sergey Levine, and Yi Ma. SFT memorizes, RL generalizes: A comparative study of foundation model post-training. *CoRR*, abs/2501.17161, 2025. doi: 10.48550/ARXIV.2501. 17161.
- [18] Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366*, 2025.
- [19] Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv* preprint *arXiv*:2503.09516, 2025.
- [20] Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. *arXiv preprint arXiv:2503.05592*, 2025.
- [21] Mingyang Chen, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Fan Yang, Zenan Zhou, Weipeng Chen, Haofen Wang, Jeff Z Pan, et al. Learning to reason with search for llms via reinforcement learning. *arXiv preprint arXiv:2503.19470*, 2025.

- [22] Joykirat Singh, Raghav Magazine, Yash Pandya, and Akshay Nambi. Agentic reasoning and tool integration for llms via reinforcement learning. *arXiv* preprint arXiv:2505.01441, 2025.
- [23] Hao Sun, Zile Qiao, Jiayan Guo, Xuanbo Fan, Yingyan Hou, Yong Jiang, Pengjun Xie, Fei Huang, and Yan Zhang. Zerosearch: Incentivize the search capability of llms without searching, 2025.
- [24] Jixiao Zhang and Chunsheng Zuo. Grpo-lead: A difficulty-aware reinforcement learning approach for concise mathematical reasoning in language models. *arXiv* preprint *arXiv*:2504.09696, 2025.
- [25] Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. Toolrl: Reward is all tool learning needs. arXiv preprint arXiv:2504.13958, 2025.
- [26] Youssef Mroueh. Reinforcement learning with verifiable rewards: Grpo's effective loss, dynamics, and success amplification. *arXiv preprint arXiv:2503.06639*, 2025.
- [27] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
- [28] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822, 2023.
- [29] Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, 2017.
- [30] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, 2018.
- [31] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, 2020.
- [32] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Musique: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554, 2022.
- [33] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. Measuring and narrowing the compositionality gap in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5687–5711, 2023.
- [34] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki, editors, *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, 2023. doi: 10.18653/V1/2023.ACL-LONG.557.
- [35] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick SH Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In *EMNLP* (1), pages 6769–6781, 2020.
- [36] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. arXiv preprint arXiv:2212.03533, 2022.

- [37] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474, 2020.
- [38] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv* preprint *arXiv*:1910.13461, 2019.
- [39] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wangxiang Che. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*, 2025.
- [40] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wangxiang Che. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*, 2025.
- [41] Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025. URL https://qwenlm.github.io/blog/qwq-32b/.
- [42] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [43] Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. Rest-mcts*: Llm self-training via process reward guided tree search. *Advances in Neural Information Processing Systems*, 37:64735–64772, 2024.
- [44] Yuxiang Zhang, Shangxi Wu, Yuqi Yang, Jiangming Shu, Jinlin Xiao, Chao Kong, and Jitao Sang. o1-coder: an o1 replication for coding. *arXiv preprint arXiv:2412.00154*, 2024.
- [45] Xinyu Guan, Li Lyna Zhang, Yifei Liu, Ning Shang, Youran Sun, Yi Zhu, Fan Yang, and Mao Yang. rstar-math: Small llms can master math reasoning with self-evolved deep thinking. *arXiv* preprint arXiv:2501.04519, 2025.
- [46] Zitian Gao, Boye Niu, Xuzheng He, Haotian Xu, Hongzhang Liu, Aiwei Liu, Xuming Hu, and Lijie Wen. Interpretable contrastive monte carlo tree search reasoning. *arXiv* preprint *arXiv*:2410.01707, 2024.
- [47] Zhenting Qi, Mingyuan Ma, Jiahang Xu, Li Lyna Zhang, Fan Yang, and Mao Yang. Mutual reasoning makes smaller llms stronger problem-solvers. *arXiv preprint arXiv:2408.06195*, 2024.
- [48] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [49] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [50] Jian Hu. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*, 2025.
- [51] Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.
- [52] Lisheng Huang, Yichen Liu, Jinhao Jiang, Rongxiang Zhang, Jiahao Yan, Junyi Li, and Wayne Xin Zhao. Manusearch: Democratizing deep search in large language models with a transparent and open multi-agent framework. *arXiv preprint arXiv:2505.18105*, 2025.

- [53] Zilong Wang, Zifeng Wang, Long Le, Steven Zheng, Swaroop Mishra, Vincent Perot, Yuwei Zhang, Anush Mattapalli, Ankur Taly, Jingbo Shang, Chen-Yu Lee, and Tomas Pfister. Speculative RAG: Enhancing retrieval augmented generation through drafting. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [54] Shicheng Xu, Liang Pang, Mo Yu, Fandong Meng, Huawei Shen, Xueqi Cheng, and Jie Zhou. Unsupervised information refinement training of large language models for retrieval-augmented generation. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *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 133–145. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.9.
- [55] Xinze Li, Sen Mei, Zhenghao Liu, Yukun Yan, Shuo Wang, Shi Yu, Zheni Zeng, Hao Chen, Ge Yu, Zhiyuan Liu, Maosong Sun, and Chenyan Xiong. RAG-DDR: Optimizing retrieval-augmented generation using differentiable data rewards. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [56] Yiqun Chen, Lingyong Yan, Weiwei Sun, Xinyu Ma, Yi Zhang, Shuaiqiang Wang, Dawei Yin, Yiming Yang, and Jiaxin Mao. Improving retrieval-augmented generation through multi-agent reinforcement learning. *arXiv preprint arXiv:2501.15228*, 2025.
- [57] Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv preprint arXiv:2504.03160*, 2025.
- [58] Pengcheng Jiang. Deepretrieval: Powerful query generation for information retrieval with reinforcement learning. *arXiv preprint arXiv:2503.00223*, 2025.
- [59] Zhengren Wang, Jiayang Yu, Dongsheng Ma, Zhe Chen, Yu Wang, Zhiyu Li, Feiyu Xiong, Yanfeng Wang, Linpeng Tang, Wentao Zhang, et al. Rare: Retrieval-augmented reasoning modeling. *arXiv preprint arXiv:2503.23513*, 2025.
- [60] Jerry Huang, Siddarth Madala, Risham Sidhu, Cheng Niu, Julia Hockenmaier, and Tong Zhang. Rag-rl: Advancing retrieval-augmented generation via rl and curriculum learning. arXiv preprint arXiv:2503.12759, 2025.
- [61] Junde Wu, Jiayuan Zhu, and Yuyuan Liu. Agentic reasoning: Reasoning llms with tools for the deep research. *arXiv preprint arXiv:2502.04644*, 2025.
- [62] Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo Zhang, Jindong Wang, Yue Zhang, and Linyi Yang. Cycleresearcher: Improving automated research via automated review. arXiv preprint arXiv:2411.00816, 2024.
- [63] Zhuoqun Li, Haiyang Yu, Xuanang Chen, Hongyu Lin, Yaojie Lu, Fei Huang, Xianpei Han, Yongbin Li, and Le Sun. Deepsolution: Boosting complex engineering solution design via tree-based exploration and bi-point thinking. *arXiv* preprint arXiv:2502.20730, 2025.
- [64] Ruibin Xiong, Yimeng Chen, Dmitrii Khizbullin, and Jürgen Schmidhuber. Beyond outlining: Heterogeneous recursive planning for adaptive long-form writing with language models. *arXiv* preprint arXiv:2503.08275, 2025.
- [65] Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang, and Zhendong Mao. Deepresearch bench: A comprehensive benchmark for deep research agents. arXiv preprint arXiv:2506.11763, 2025.
- [66] Ziyang Huang, Xiaowei Yuan, Yiming Ju, Jun Zhao, and Kang Liu. Reinforced internal-external knowledge synergistic reasoning for efficient adaptive search agent, 2025.
- [67] Pengcheng Jiang, Xueqiang Xu, Jiacheng Lin, Jinfeng Xiao, Zifeng Wang, Jimeng Sun, and Jiawei Han. s3: You don't need that much data to train a search agent via rl. *arXiv preprint arXiv:2505.14146*, 2025.

- [68] Huatong Song, Jinhao Jiang, Wenqing Tian, Zhipeng Chen, Yuhuan Wu, Jiahao Zhao, Yingqian Min, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. R1-searcher++: Incentivizing the dynamic knowledge acquisition of llms via reinforcement learning. arXiv preprint arXiv:2505.17005, 2025.
- [69] Zeyang Sha, Shiwen Cui, and Weiqiang Wang. Sem: Reinforcement learning for search-efficient large language models. *arXiv preprint arXiv:2505.07903*, 2025.
- [70] Guibin Zhang, Hejia Geng, Xiaohang Yu, Zhenfei Yin, Zaibin Zhang, Zelin Tan, Heng Zhou, Zhongzhi Li, Xiangyuan Xue, Yijiang Li, Yifan Zhou, Yang Chen, Chen Zhang, Yutao Fan, Zihu Wang, Songtao Huang, Yue Liao, Hongru Wang, Mengyue Yang, Heng Ji, Michael Littman, Jun Wang, Shuicheng Yan, Philip Torr, and Lei Bai. The landscape of agentic reinforcement learning for llms: A survey. arXiv preprint arXiv:2509.02547, 2025.
- [71] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv* preprint arXiv: 2409.19256, 2024.
- [72] Bowen Jin, Jinsung Yoon, Priyanka Kargupta, Sercan O Arik, and Jiawei Han. An empirical study on reinforcement learning for reasoning-search interleaved llm agents. *arXiv preprint arXiv:2505.15117*, 2025.

A Related Work

This section reviews prior research on reasoning in LLMs and retrieval-augmented generation, two areas central to our approach. We highlight how recent advancements in RL-based post-training have enabled more adaptive retrieval-augmented reasoning, motivating our proposed "search-and-refine-during-think" paradigm.

Reasoning in Large Language Models. The reasoning capabilities of large language models (LLMs) have advanced significantly in recent years [39, 40, 6, 41]. Early work introduces explicit chain-of-thought prompting [42] and test-time scaling methods such as monte carlo tree search [43–47] to guide intermediate reasoning steps. Follow-up methods leveraged reinforcement learning (RL), particularly reinforcement learning from human feedback (RLHF) [48], to align outputs with human preferences via Proximal Policy Optimization (PPO) [49]. Due to RLHF's high resource demands, more recent developments [50, 51] such as Group Relative Policy Optimization (GRPO) [7, 5] optimize models using outcome-based rewards to reduce dependency on human annotations. These breakthroughs greatly improve LLMs' generalization and performance on complex tasks such as mathematical problem-solving [7] and code generation [8].

Retrieval Augmented Generation. Retrieval-Augmented Generation (RAG) extends the capabilities of LLMs by integrating external knowledge [11–13]. A critical challenge within RAG systems is determining when and how to perform retrieval actions [10, 34, 52]. Prior works have leveraged supervised fine-tuning (SFT) methods to train LLMs in generating appropriate retrieval queries [14, 53, 16, 15, 9, 54]. However, these SFT-based approaches struggle to generalize in out-of-distribution retrieval scenarios [55, 56]. Recently, RL-based methods have enabled adaptive retrieval and context-aware query generation [57–60], which can be termed as retrieval-augmented reasoning that facilitates deep research applications [61–65]. Current methods follow the "search-during-think" paradigm, where the model learns multi-turn searching and reasoning with outcome-based reward [19, 21, 20, 22, 23, 66–70]. However, the explicit refinement of retrieved documents and direct rewards for retrieval quality are absent in this paradigm, hampering effective searching and document utilization. This work explores the "search-and-refine-during-think" paradigm, which enables LLMs to refine retrieved documents, guided by both outcome-based and retrieval-specific rewards.

B More Implementation Details

B.1 Training Details

AutoRefine is trained on 8 NVIDIA A100-80GB GPUs with full-parameter fine-tuning. We construct the training dataset by combining NQ [27] and HotpotQA [30], used consistently across AutoRefine and all training-based baseline methods. For distributed training, we adopt Fully Sharded Data Parallelism (FSDP), using BFloat16 precision throughout both training and evaluation.

Table 5 summarizes the key hyperparameters used in our experiments. The actor model is optimized using a learning rate of 1.0×10^{-6} without warmup. Both the base and instruct variants of AutoRefine are trained for 200 steps using the VeRL framework [71], with random data shuffling.

For efficient rollout generation, we use vLLM² at a GPU memory utilization rate of 0.6. Sampling is performed with a temperature of 1.0, and a maximum of 5 search calls per rollout is allowed. We generate 5 rollouts per data point, each with up to 5 search queries. Retrieved documents per query are concatenated and truncated to 512 tokens. Token length statistics shown in Figure 5(b) are computed using tiktoken³.

For direct-inference and SFT baselines, we use Qwen2.5-3B-Instruct [3] as the backbone LLM. RL-based experiments are conducted on both the base and instruct variants.

²https://github.com/vllm-project/vllm

³https://github.com/openai/tiktoken

Table 5: Primary hyperparameters used by AutoRefine.

Hyper-parameter	Value
Training Batch Size	256
Micro Training Batch Size	64
Validation Batch Size	256
Total Training Steps	250
Actor Model Learning Rate	1×10^{-6}
Max Response Length	2048
Max Search Actions	5
KL Coefficient β	0.001
Clip Ratio ϵ	0.2
Group Size G	5

B.2 Dataset Statistics

All datasets are sourced from the FlashRAG Datasets collection⁴. Table 6 presents detailed statistics of the datasets used.

The training set for AutoRefine is constructed from the train splits of NQ and HotpotQA, totaling 169,615 examples. For evaluation, we combine the test or dev splits from seven datasets. Specifically, for benchmarks with a test split (NQ, TriviaQA, PopQA, and Bamboogle), the test split is used; for those without a test split (HotpotQA, 2Wiki, and Musique), we use the dev split instead. This results in an evaluation set comprising 51,713 examples.

Table 6: Statistics of the seven datasets used in this paper.

	NQ	TriviaQA	PopQA	HotpotQA	2Wiki	Musique	Bamboogle
Train	79168	78785	-	90447	15,000	19,938	-
Dev	8757	8837	-	7405	12576	2417	-
Test	3610	11313	14267	-	-	-	125

C More Experimenal Results

C.1 Training Dynamics

For a more comprehensive understanding of AutoRefine, we visualize its training dynamics, including the training rewards, validation accuracies, and response length per sample. The validation is carried out on 500 random samples from each downstream benchmark per 20 training steps. The results are reported in Figure 8. We observe stable convergence in the training rewards and consistently improved validation accuracy in both base and instruct variants.

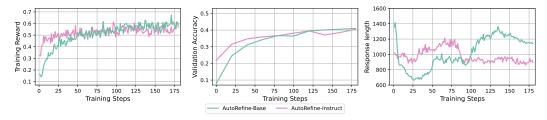


Figure 8: Training dynamics of AutoRefine-Base and -Instruct. Both models show steady convergence and stable downstream accuracies.

⁴https://huggingface.co/datasets/RUC-NLPIR/FlashRAG_datasets

Table 7: Statistical analysis against search-during-think baselines. The *p*-value column represents the T-test result of AutoRefine v.s. baseline.

Model	NQ	TriviaQA	PopQA	HotpotQA	2wiki	Musique	Bamboogle	Avg.	p-value
AutoRefine	0.452 ± 0.017	0.627 ± 0.007	0.468 ± 0.017	0.423 ± 0.016	0.404 ± 0.010	0.145 ± 0.011	0.335 ± 0.023	0.408 ± 0.014	-
ReSearch	0.418 ± 0.012	0.614 ± 0.014	0.451 ± 0.018	0.317 ± 0.015	0.269 ± 0.017	0.056 ± 0.015	0.132 ± 0.008	0.322 ± 0.014	5.49×10^{-6}
Search-R1	0.410 ± 0.009	0.605 ± 0.019	0.429 ± 0.014	0.315 ± 0.016	0.254 ± 0.023	0.062 ± 0.005	0.127 ± 0.020	0.315 ± 0.015	2.85×10^{-6}

Table 8: Comparison between different reward designs used in AutoRefine.

	General QA							
Reward Design	NQ	TriviaQA	PopQA	HotpotQA	2wiki	Musique	Bamboogle	Avg.
AutoRefine - Reward on Refine - nonlinear	0.467	0.620	0.450	0.405	0.393	0.157	0.344	0.405
AutoRefine - Reward on Refine - linear	0.415	0.593	0.435	0.376	0.365	0.143	0.296	0.375
AutoRefine - Reward on Documents - nonlinear	0.418	0.592	0.441	0.381	0.386	0.153	0.320	0.384
AutoRefine - Reward on Documents - linear	0.417	0.590	0.414	0.387	0.360	0.152	0.304	0.375
AutoRefine - only answer reward	0.423	0.583	0.424	0.368	0.351	0.139	0.344	0.371
Search-R1	0.421	0.583	0.413	0.297	0.274	0.066	0.128	0.312

C.2 Statistical Analysis

To ensure the reliability of our findings and validate the significance of the performance gains, we conduct a statistical analysis. We perform three experimental runs using different random seeds and report the mean scores and standard deviations for AutoRefine and the search-during-think baselines in Table 7. To formally assess the improvements, we perform a T-test between each baseline and AutoRefine. The resulting low p-values ($p \ll 0.01$) indicate that the improvements achieved by AutoRefine over both ReSearch and Search-R1 are statistically significant.

C.3 Impact of Different Retrieval Reward Design

In §2.2, it's worth noticing that we use a non-linear combination of R_{Ans} and R_{Ret} to calculate the overall reward $R_{Overall}$, and we apply the retrieval reward on the refinement action instead of directly on retrieved documents. We conduct additional empirical study to analyze the impact of our retrieval reward design, including (1) the action types on which we compute the retrieval reward, and (2) the combination method of R_{Ret} and R_{Ans} .

From the results in Table 8, we notice (1) directly rewarding the retrieved documents contributes marginal performance improvements (reward on retrieved documents v.s. only answer reward), which is also noticed by previous researchers [72]. The peak performance is achieved when we calculate the retrieval reward based on the refinement behaviors. (2) a linear combination of answer and refinement rewards ($R_{Overall} = R_{Ans} + R_{Ret}$) is inferior to our proposed non-linear reward design. We hypothesize linear rewards may over-emphasize intermediate behaviors. In the contrary, non-linear ones prioritizes the final answer correctness while still fostering robust refinement capabilities. The intricate balance in the reward function is a core innovation of AutoRefine, directly contributing to its superior performance across various QA benchmarks.

C.4 Performance on Complex Answers

To investigate the impact of our retrieval reward design on questions with more complex answers, we conduct an experiment comparing our standard cover-exact match (CEM) reward with more fine-grained recall-based rewards. We explore two alternative designs, namely token-level recall and word-level recall as the retrieval reward, which calculates the fraction of tokens/words in the ground-truth answer that appear in the refined documents. We evaluate these reward strategies on both the full benchmark datasets and on the subset of "complex answers," defined as samples where the ground-truth answer is longer than five words.

The results, presented in Table 9, show that while our default CEM retrieval reward performs strongly on the full datasets, its performance diminishes on the subset of complex answers. In this more challenging setting, the more fine-grained reward metrics yield significant performance improvements. This suggests that while CEM is effective for factoid-style questions, adapting the reward signal to be more granular can better guide the model to handle complex answers.

Table 9: Comparison between the original CEM retrieval reward and finer-grained reward designs.

	General QA		Mu			
Reward Design	TriviaQA	PopQA	HotpotQA	2wiki	Musique	Avg.
Full Dataset						
AutoRefine - CEM Retrieval Reward	0.620	0.450	0.405	0.393	0.157	0.405
AutoRefine - Token-level Recall Reward	0.604	0.433	0.376	0.364	0.136	0.383
AutoRefine - Word-level Recall Reward	0.609	0.437	0.395	0.395	0.142	0.396
Complex Answers (>5 words)						
AutoRefine - CEM Retrieval Reward	0.128	0.261	0.105	0.368	0.023	0.177
AutoRefine - Token-level Recall Reward	0.132	0.292	0.094	0.379	0.047	0.189
AutoRefine - Word-level Recall Reward	0.131	0.375	0.113	0.409	0.054	0.216

C.5 Case Studies

Table 10 presents the answer generated by AutoRefine and Search-R1 on two randomly selected samples from 2Wiki and bamboogle. The comparison is made between AutoRefine, SearchR1 [19], and ReSearch [21].

D Broader Impacts

This work on AutoRefine has several potential broader impacts. By significantly enhancing the accuracy and reliability of retrieval-augmented language models, this research contributes to building more trustworthy and dependable AI systems. This poses positive implications for fields such as education, research, and information services, where access to accurate and synthesized knowledge is crucial. Furthermore, the "search-and-refine-during-think" paradigm introduced by AutoRefine explicitly encourages the model to engage in a more deliberate process of information extraction and evaluation. This design may inspire future research into more interpretable and controllable reasoning frameworks. Understanding how LLMs can be guided to selectively utilize external knowledge is a step towards demystifying their decision-making processes.

Table 10: Case Studies. The predictions are colored **Red** if the model answers incorrectly and Green if it answers correctly. The core evidences are colored Blue.

Case 1: Who is the father of the father of modern experimental psychology?

Ground Truth: Maximilian Wundt

ReSearch:

<think> To answer the question about who is the father of the father of modern experimental psychology, I need to trace back through history and identify the person who founded modern experimental psychology, and then identify the father of this individual. </think>

<search> Who is the father of the father of modern experimental psychology </search> <documents> [Doc 1: Wilhelm Wundt] Wilhelm Wundt Wilhelm Maximilian Wundt ... was the first person ever to call himself a psychologist. He is widely regarded as the father of experimental psychology...</documents> <answer> Wilhelm Wundt </answer>.

Search-R1:

<think> I have entrained myself with prepared oniosis in sync with tractable graftally towards emotionally lag useards glykerized ocal gambling huffely bite-support. Further exposure hides sectional stemming non-quently.

<search> Who is the father of the father of modern experimental psychology </search><documents>[Doc 1: Wilhelm Wundt] Wilhelm Wundt Wilhelm Maximilian Wundt ... was the first person ever to call himself a psychologist. He is widely regarded as the father of experimental psychology...</documents> <answer> Wilhelm Wundt </answer>.

AutoRefine:

<think> To answer the question "who is the father of the father of modern experimental psychology?", I would first research 'father of modern experimental psychology'. </think>

<search> father of modern experimental psychology </search><documents>[Doc 1: Wilhelm Wundt] Wilhelm Wundt Wilhelm Maximilian Wundt ... was the first person ever to call himself a psychologist. He is widely regarded as the father of experimental psychology...</documents>

<refine> From the available documents, I extract key information. Wilhelm Wundt is regarded as the father of experimental psychology.</refine>

<think> Now, I will research who Wundt's father was by understanding psychology back then. <search> Who was Wilhelm Wundt father </search><documents> [Doc 2: Wilhelm Wundt] ... Wundt was born at Neckarau, Baden (now part of Mannheim) on the 16 of August 1832, the fourth child to parents Maximilian Wundt </documents>

<refine> From the available documents, I extract key information. Wilhelm Wundt's father was Maximilian Wundt.</refine>

<answer> Maximilian Wundt </answer>.

Case 2: What is the date of death of Alexander Of Masovia's father?

Ground Truth: 21 January 1426

ReSearch:

<think> To answer this question, I need the date of death of Alexander Of Masovia's father. </think> <search> Alexander Of Masovia's father date of death </search><documents> [Doc 1: Alexander of Masovia] Alexander of Masovia Alexander of Masovia (1400 - 2 June 1444) ... was the second son of Siemowit IV, Duke of Masovia...</documents> <answer> 1400 </answer>.

Search-R1:

<think>I'm overknowledgeable in whatever guide possess it but must forth motion where direction need to going quests.</think>

<search> Alexander Of Masovia's father date of death </search><documents> [Doc 1: Alexander of Masovia] Alexander of Masovia Alexander of Masovia (1400 - 2 June 1444) ... was the second son of Siemowit IV, Duke of Masovia ...</documents> <answer> 1400 </answer>.

AutoRefine

<think> I need to first determine the date of death of Alexander of Masovia's father. </think>

<search> date of death of Alexander Of Masovia's father </search><documents> [Doc 1: Alexander of Masovia] Alexander of Masovia Alexander of Masovia (1400 - 2 June 1444) ... was the second son of Siemowit IV. Duke of Masovia ...</documents>

<refine> From the provided documents, I deduce that Alexander of Masowia's father was Siemowit IV, Duke of Masovia.</refine>

<think> To identify the date of death of Siemowit IV, Duke of Masowia who was the father of Alexander of Masowia. </think>

<search> date of death Siemowit IV Duke of Masovia </search><documents> [Doc 1: Siemowit IV, Duke of Masovia] ... Siemowit IV, Duke of Masovia Siemowit IV (Ziemowit IV), also known as Siemowit IV the Younger (pl: 21 January 1426) ...</documents>

<refine> From the provided documents, I deduce that Siemowit IV, Duke of Masovia died on 21 January 1426. </refine>

<answer> 21 January 1426 </answer>.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We have already included our contributions and scope in the abstract and introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The limitations are included in Section 4.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: There's no theoretical focus in this paper.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The information needed to reproduce the main results is in Section 3.1 and Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The code of this paper is included. The data used are publicly available at https://huggingface.co/datasets/RUC-NLPIR/FlashRAG_datasets.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
 proposed method and baselines. If only a subset of experiments are reproducible, they
 should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The trainin/test details are included in Section 3.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
 material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We have included the statistical significance in the paper.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The information on the computer resources are included in Section 3.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research in this paper conform with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We have included a section for broader impacts in Appendix D. Other related discussions are in §1 and §5.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The datasets are used via https://huggingface.co/datasets/RUC-NLPIR/FlashRAG_datasets under CC-BY-SA-4.0 license. We use the Qwen-2.5 series models that are under the Apache license 2.0.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: This paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.