

SmartSearch: Process Reward-Guided Query Refinement for Search Agents

Tongyu Wen
wentongyu@ruc.edu.cn
Renmin University of China
Beijing, China

Guanting Dong
dongguanting@ruc.edu.cn
Renmin University of China
Beijing, China

Zhicheng Dou
dou@ruc.edu.cn
Renmin University of China
Beijing, China

Abstract

Large language model (LLM)-based search agents have proven promising for addressing knowledge-intensive problems by incorporating information retrieval capabilities. Existing works largely focus on optimizing the reasoning paradigms of search agents, yet the quality of intermediate search queries during reasoning remains overlooked. As a result, the generated queries often remain inaccurate, leading to unexpected retrieval results and ultimately limiting search agents' overall effectiveness. To mitigate this issue, we introduce **SmartSearch**, a framework built upon two key mechanisms: (1) **Process rewards**, which provide fine-grained supervision for the quality of each intermediate search query through Dual-Level Credit Assessment. (2) **Query refinement**, which promotes the optimization of query generation by selectively refining low-quality search queries and regenerating subsequent search rounds based on these refinements. To enable the search agent to progressively internalize the ability to improve query quality under the guidance of process rewards, we design a three-stage curriculum learning framework. This framework guides the agent through a progression from imitation, to alignment, and ultimately to generalization. Experimental results show that SmartSearch consistently surpasses existing baselines, and additional quantitative analyses further confirm its significant gains in both search efficiency and query quality. The code is available at <https://github.com/MYVAE/SmartSearch>.

CCS Concepts

• **Information systems** → **Information retrieval**; **Language models**; **Question answering**.

Keywords

Search Agent, Information Retrieval, Large Language Models, Process Reward, Query Refinement

ACM Reference Format:

Tongyu Wen, Guanting Dong, and Zhicheng Dou. 2026. SmartSearch: Process Reward-Guided Query Refinement for Search Agents. In *Proceedings of the 49th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '26)*, July 20–24, 2026, Melbourne | Naarm, Australia.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '26, Melbourne | Naarm, Australia

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/2026/06
<https://doi.org/XXXXXXX.XXXXXXX>

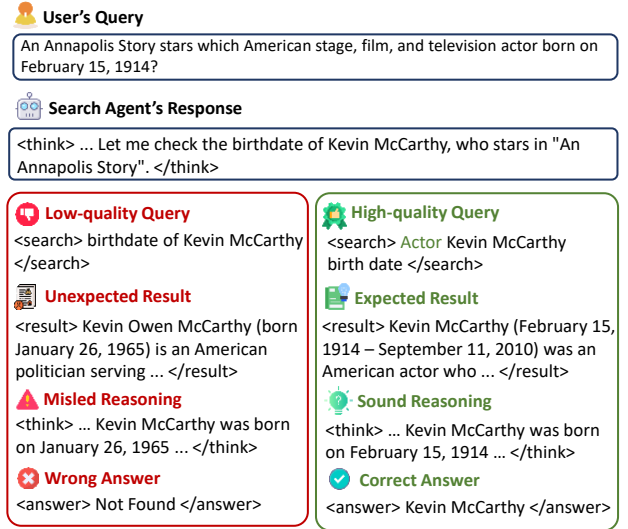


Figure 1: An example from ASearcher [14] dataset demonstrating how low-quality intermediate search queries lead to unexpected retrieval results and derail the entire trajectory.

Australia. ACM, New York, NY, USA, 16 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Large language models (LLMs) have shown strong performance across a variety of tasks [1, 2, 6, 45, 46], including translation [58, 64], summarization [41, 66], and question answering [23, 38]. However, challenges remain, particularly with issues like hallucinations [17, 35] and the absence of recent or field-specific knowledge, which may result in inaccurate or outdated answers. Retrieval-augmented generation (RAG) [3, 15, 19] has been introduced to address these challenges by incorporating external knowledge to complement the model's internal knowledge [36, 54]. However, static RAG faces limitations in its ability to handle more complex, dynamic, and deep exploration tasks.

Recently, LLM-based search agents have proven to be a promising method [4, 21, 28, 29, 39, 40]. These agents can autonomously and iteratively invoke external search tools, thereby addressing more challenging knowledge-intensive problems that demand adaptive retrieval and in-depth reasoning. Current research on search agents has made considerable progress in optimizing the reasoning paradigms of search agents through methods like prompt engineering [28] and fine-tuning [4, 21, 39, 40]. However, they often overlook

the quality of intermediate search queries during reasoning, yet low-quality queries can lead to unexpected retrieval results or even derail the entire trajectory. Figure 1 illustrates how minor inaccuracies in an intermediate search query (e.g., omitting ‘actor’) can lead a search agent to retrieve and accept unexpected information, ultimately resulting in an incorrect answer. This highlights the critical role that search query quality plays in the deep information-seeking process. Some studies [9, 52, 59, 67] have attempted to incorporate process rewards into search agent training. However, they tend to focus more on shaping better reasoning behavior rather than improving the quality of intermediate search queries, and existing efforts [52] on intermediate search queries remain preliminary and ineffective. Furthermore, research [20, 42] has shown that existing training paradigms often prioritize information utilization, persistently neglecting the optimization of retrieval patterns. This undoubtedly impedes the search agent’s ability to achieve deep and reliable information retrieval, thereby compromising its overall effectiveness. Such issues highlight the need for methods that specifically focus on optimizing query quality during training.

In this work, we present **SmartSearch**, a framework that optimizes search query quality through the guidance of process rewards, thereby enhancing the deep information-seeking capabilities of search agents. Specifically, SmartSearch incorporates two key mechanisms: **(1) Process rewards**: To provide fine-grained supervision for the quality of each search query, we introduce Dual-Level Credit Assessment, which comprises two complementary components. The first one is a rule-based assessment for query novelty, which detects redundancy by checking whether the retrieved documents contain excessive overlap with previous rounds. The second one is a model-based evaluation for query usefulness, which judges whether the query intent is necessary and whether the retrieved results provide the expected answer. This mechanism outputs both numerical scores and textual feedback, which serve as guidance for subsequent query refinement. **(2) Query refinement**: To further promote the optimization of query generation during training, the agent first generates a complete search trajectory, then identifies low-quality search rounds according to the numerical scores from the process rewards. Subsequently, we employ a model to refine those queries under the textual guidance provided by the process rewards, after which the search agent continues generating from the refined queries. To improve the efficiency of query assessment and refinement, a smaller model is trained for scoring and refinement, reducing computational cost while maintaining effectiveness.

Building on the foundation of the two mechanisms, we introduce a three-stage curriculum learning framework. The framework guides the search agent through a progression from imitation and alignment to generalization, enabling it to progressively internalize the ability to enhance query quality under the guidance of process rewards. **(1) Query Quality Screened Imitation Learning**: The initial stage leverages Supervised Fine-Tuning (SFT) to guide the search agent during its early learning of information retrieval and utilization. The training data is filtered based on both final answer correctness and query quality measured by the process rewards. It ensures the model to learn from trajectories that not only lead to correct answers but also maintain high-quality search processes. **(2) Query Generation Alignment**: In this stage, the search agent cultivates advanced query generation capabilities through Direct

Preference Optimization (DPO). We employ the query refinement mechanism to generate comparative data, with process rewards and outcome rewards jointly defining which trajectories are of higher quality. **(3) Query-Aware Policy Optimization**: The final stage utilizes Reinforcement Learning (RL) to further strengthen its integrated capabilities of information retrieval and utilization. During the rollout phase, the query refinement mechanism is employed, with the process rewards incorporated into the reward function.

To thoroughly assess the capabilities of SmartSearch, we perform experiments on four challenging knowledge-intensive tasks and two web exploration tasks. Experimental results indicate that SmartSearch consistently surpasses all baselines in overall performance and exhibits strong generalization to open-web settings. Additionally, we perform a range of ablation studies and quantitative analyses to comprehensively validate SmartSearch’s effectiveness. Our findings highlight the critical contribution of our two key mechanisms and three curriculum learning stages, as well as their superiority in terms of search efficiency, search query quality, and other dimensions.

To summarize, the primary contributions of this study include:

- (1) We present a pioneering focus that optimizes the quality of intermediate search queries through process reward guidance, thereby improving the information-seeking ability of search agents.
- (2) We propose SmartSearch, a framework that incorporates two key mechanisms: process rewards and query refinement, to enable process reward-guided search refinement.
- (3) We design a three-stage, query-oriented curriculum learning framework that guides the agent through a progression from imitation and alignment to generalization, progressively internalizing the ability to improve query quality.
- (4) Experiments across six challenging benchmarks demonstrate that SmartSearch consistently surpasses existing baselines, and further quantitative analyses confirm significant improvements in both search efficiency and query quality.

2 Related Works

2.1 LLM-based Search Agents

LLMs have demonstrated strong performance across various tasks [1, 2, 6, 45, 46], yet challenges like hallucinations [17, 35] and static parametric knowledge remain. Nowadays, LLM-based search agents have emerged as a promising solution [4, 21, 28, 29, 39, 40]. This advanced paradigm enables models to autonomously and iteratively invoke external tools, effectively tackling challenging knowledge-intensive problems. Research on search agents has progressed through methods including prompt engineering and fine-tuning. Early prompt-based methods [25, 28, 31, 48] focused on carefully designed prompts and structured workflows to steer the agent’s behavior. However, these methods don’t fundamentally enhance the model’s underlying capabilities, leading many studies to shift towards fine-tuning-based approaches. A prominent line of work has demonstrated that SFT [12, 13, 18, 30] on expert trajectories enables agents to learn through imitation and yields promising performance. Building upon this foundation, recent studies [4, 10, 21, 39, 40] have employed RL to further advance search agent capabilities. However, existing methods tend to overlook intermediate search query quality, which can lead to unexpected retrieval results or even derail

the entire trajectory. Moreover, research [20] indicates that current training paradigms tend to prioritize information utilization, which can lead to stagnation in information retrieval abilities. Thus, we present a framework designed to optimize the quality of intermediate search queries under the guidance of process rewards, thereby enhancing the overall performance of search agents.

2.2 Process Rewards in RL

Recent advancements in RL have achieved significant success in large reasoning models [7, 43, 44], and have also demonstrated effectiveness in enhancing the performance of LLM-based search agents [4, 10, 21, 28, 29, 39, 40]. However, reward signals based solely on final outcomes often result in sparse feedback in multi-round search tasks, providing insufficient guidance for intermediate steps and leading to unstable and inefficient policy optimization [65]. To overcome this limitation, recent studies [9, 52, 59, 67] have explored the use of process-based rewards. Some approaches employ Monte Carlo Tree Search to estimate intermediate actions' value [26, 67], while others rely on a corpus of annotated golden steps or intermediate information to compute rewards based on alignment with this reference [49, 51, 52, 63, 68]. Still others leverage external reward models to provide fine-grained evaluation for each step [9, 56, 57, 59, 60]. These approaches have proven effective in enhancing the effectiveness and stability of RL training. Yet, most of these approaches tend to focus primarily on the quality of the reasoning process rather than the quality of intermediate search queries [9, 60], with existing efforts on intermediate search queries remaining preliminary and ineffective [52]. In this context, our process rewards mechanism provides fine-grained supervision for query quality through Dual-Level Credit Assessment, playing a central role in the query-oriented training framework.

3 Preliminaries

3.1 Task Formulation

We adopt ReAct [62] as the framework for the search agent. Given a user query q , the search agent, guided by an LLM policy π_θ , interacts with an external search tool through several iterations of Thought-Action-Observation to gather information and ultimately generate an answer. Specifically, during each iteration, the search agent starts by engaging in thinking to generate a "Thought" according to the existing context. It then produces the next "Action", which involves querying the search tool. The agent subsequently waits for the environment to return the "Observation", consisting of the Top-K retrieved document fragments for the search query. The iteration concludes when the search agent has gathered sufficient information required to address the user's question and selects the "final answer" as the action. A complete trajectory over T iterations is denoted as:

$$H_T = (q, \tau_0, a_0, o_0, \dots, \tau_i, a_i, o_i, \dots, \tau_T, a_T). \quad (1)$$

Here, τ_i , a_i , and o_i correspond to the Thought, Action, and Observation of the i -th iteration. In iteration t , the LLM policy $\pi_\theta(a, t|H_{t-1})$ produces the thought τ_t and action a_t , which is conditioned on the entire history of prior context H_{t-1} .

3.2 Agentic Reinforcement Learning

Policy Optimization. In the context of Agentic RL [65], Group Relative Policy Optimization (GRPO) [37] is typically employed for policy optimization. In our approach, we also employ GRPO during the Query Aware Policy Optimization stage, with a specific focus on the augmentation of the rollout and reward modules to optimize the quality of intermediate search queries. Specifically, GRPO optimizes the policy model through maximization of the objective function below:

$$J_{\text{GRPO}}(\theta) = \mathbb{E}_{(q,a) \sim \mathcal{D}, \{o_i\} \sim \pi_{\theta_{\text{old}}}(\cdot|q)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \min \left(r_t(\theta) \hat{A}_i, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i \right) - \beta D_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}}) \right]. \quad (2)$$

In this formulation, for each input pair (q, a) drawn from the dataset \mathcal{D} , G trajectories $\{o_i\}_{i=1}^G$ are generated from the old policy $\pi_{\theta_{\text{old}}}(\cdot|q)$. The importance weight $r_t(\theta)$ is defined as:

$$r_t(\theta) = \frac{\pi_\theta(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t} | q, o_{i,<t})}. \quad (3)$$

The normalized advantage score \hat{A}_i is denoted as:

$$\hat{A}_i = \frac{r_i - \text{mean}(\{r_j\}_{j=1}^G)}{\text{std}(\{r_j\}_{j=1}^G)}. \quad (4)$$

Here, r_i denotes the scalar reward for the i -th rollout. Furthermore, agentic RL typically masks observations originating from the external environment during loss computation, thereby preventing unstable training.

Reward Design. As discussed earlier, in agentic RL, each rollout corresponds to a scalar reward r . Prior research [4, 21, 39, 40] predominantly relies on combining two key types of rewards: the outcome reward r_{outcome} , reflecting the trajectory's answer correctness, and the format reward r_{format} , assessing the trajectory's structural correctness. These rewards are typically weighted and combined using a simple hyperparameter λ as follows:

$$r = r_{\text{outcome}} + \lambda \cdot r_{\text{format}}. \quad (5)$$

In some recent works [9, 59, 60], process rewards have been incorporated into the reward function to provide fine-grained feedback on intermediate steps. The reward function is then extended to include the process rewards, with a composite reward incorporating both the outcome reward and the process rewards, while the format reward is weighted by a hyperparameter:

$$r = r_{\text{composite}} + \lambda \cdot r_{\text{format}}. \quad (6)$$

Here, $r_{\text{composite}}$ is computed as the aggregation of multiple step-wise process rewards and the final outcome reward r_{outcome} :

$$r_{\text{composite}} = f(r_1^{\text{process}}, r_2^{\text{process}}, \dots, r_n^{\text{process}}, r_{\text{outcome}}), \quad (7)$$

where n represents the total steps in the trajectory, and r_i^{process} denotes the process reward for the i -th step. The aggregation function f combines these individual rewards, and its specific form may vary across different works.

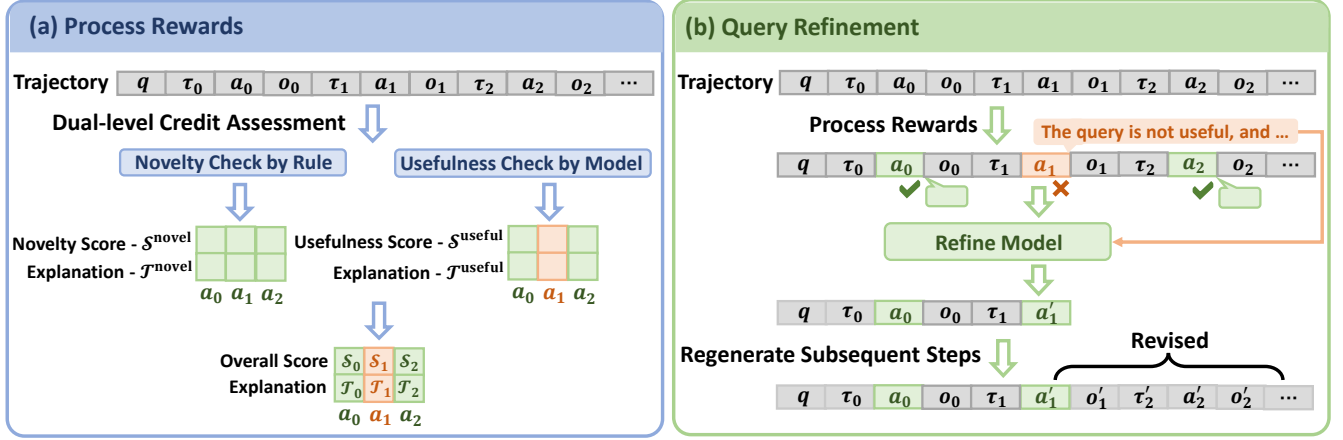


Figure 2: An overview of the two key mechanisms in SmartSearch: the process rewards (a) and the query refinement (b).

4 Our Method

4.1 Overview

We propose SmartSearch, a framework that enhances search agent performance by optimizing the quality of intermediate search queries through process reward guidance. As illustrated in Figure 2, SmartSearch incorporates two key mechanisms: (1) **Process rewards** (§4.2), which provide fine-grained supervision for the quality of each query through Dual-Level Credit Assessment. (2) **Query refinement** (§4.3), which promotes the optimization of query generation by selectively refining low-quality queries and regenerating subsequent search rounds based on these refinements. To further internalize the ability to improve query quality, we propose a three-stage curriculum learning framework (§4.4) built upon these mechanisms. As shown in Figure 3, it comprises **Query Quality Screened Imitation Learning**, **Query Generation Alignment**, and **Query Aware Policy Optimization**. Below, we will first introduce the two key mechanisms, followed by a detailed description of the three-stage curriculum learning framework.

4.2 Process Reward for Assessing Query Quality

In this section, we introduce the process rewards mechanism to assess the quality of each query, providing both numerical scores and textual feedback. These outputs guide the subsequent query refinement, and play a key role within the three-stage curriculum learning framework by selecting trajectories with high-quality search processes and providing finer-grained supervision signals, which will be introduced later.

Design Principles. Our assessment of search query quality is guided by a comprehensive set of three fundamental principles:

- **Query Novelty:** The query should avoid redundancy with previous queries and introduce novel information.
- **Intent Necessity:** The query’s search intent must be necessary for progressing toward the final answer.
- **Retrieval Relevance:** The retrieved documents should align with the search intent, effectively containing the expected information or answer.

These principles are well-motivated and collectively capture the essential aspects of a high-quality query, while also being readily applicable via either rule-based checks or simple model judgments.

Dual-Level Credit Assessment. We operationalize these principles through Dual-Level Credit Assessment, which consists of two complementary components.

(1) Rule-based Evaluation: The first is a rule-based evaluation for query novelty, which identifies redundant queries by measuring the document overlap between the current and previous search rounds. Formally, for the t -th step, the novelty score S_t^{novel} and its corresponding textual explanation $\mathcal{T}_t^{\text{novel}}$ are defined as:

$$(S_t^{\text{novel}}, \mathcal{T}_t^{\text{novel}}) = \begin{cases} (0, \text{the query is redundant}), & \text{if } O_t > K, \\ (1, \text{the query is novel}), & \text{if } O_t \leq K. \end{cases} \quad (8)$$

Here, K is a threshold hyperparameter, and O_t represents the number of documents retrieved at step t that share the same content with those retrieved in any previous step, defined as:

$$O_t = \sum_{i=1}^n \mathbb{I}(D_i^t \in \bigcup_{s=0}^{t-1} \bigcup_{j=1}^n D_j^s), \quad (9)$$

where D_i^t refers to the i -th document retrieved at step t , and $\mathbb{I}(\cdot)$ is the indicator function.

(2) Model-based Evaluation: The second component is model-based evaluation for query usefulness, which assesses the necessity of the query intent and checks whether the retrieved results provide the expected answer. For the t -th step, the evaluation score S_t^{useful} and its corresponding textual explanation $\mathcal{T}_t^{\text{useful}}$ are defined as:

$$S_t^{\text{useful}}, \mathcal{T}_t^{\text{useful}} = \text{LLM}_{\text{eval}}(q, a, H_t), \quad (10)$$

where LLM_{eval} is the model used for evaluation, q is the user’s query, a denotes the golden answer, and H_t indicates the trajectory up to step t . The score S_t^{useful} is set to 1 if the query meets the criteria, and 0 otherwise, while the explanation $\mathcal{T}_t^{\text{useful}}$ is directly parsed from the model’s output. To enhance efficiency, we employ a smaller model fine-tuned via SFT for both scoring and the subsequent query refinement task. Specifically, we input task-specific prompts into

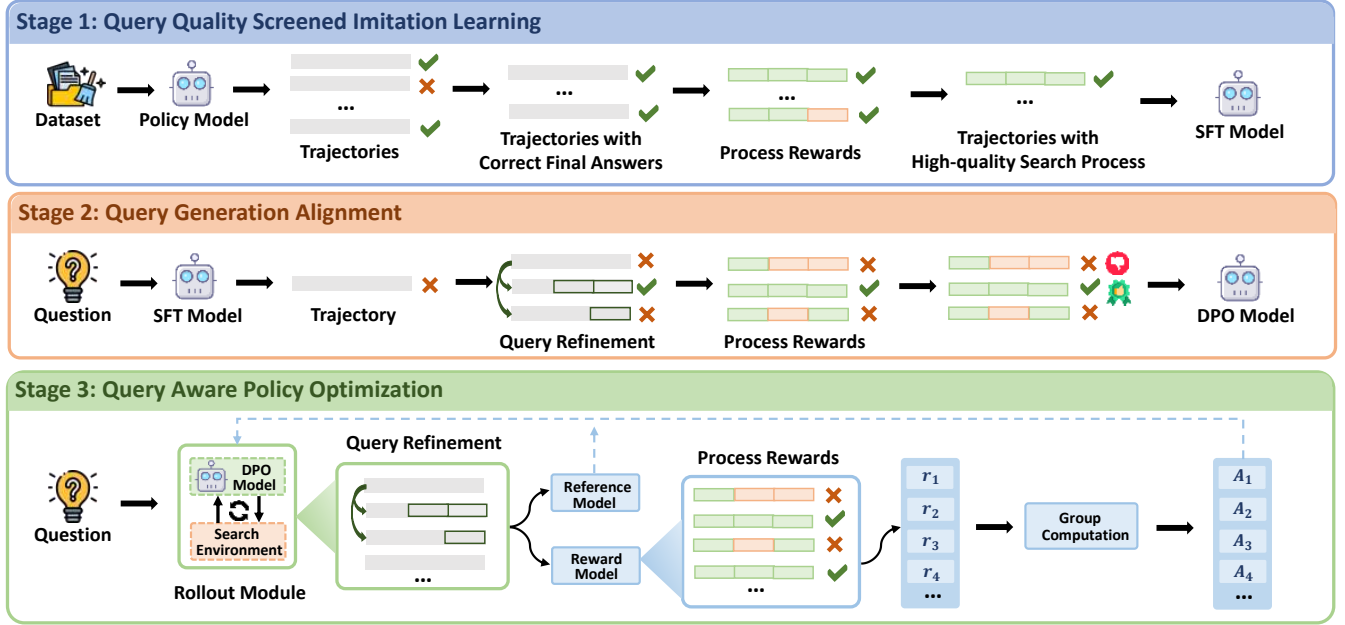


Figure 3: The overall framework of query-oriented three-stage curriculum learning, including Query Quality Screened Imitation Learning, Query Generation Alignment, and Query Generation Alignment.

a more powerful teacher model and use its outputs as annotation labels. The smaller model is then trained on these prompt-output pairs, enabling it to achieve effective performance at a reduced computational cost. More details about the model will be introduced in Section 5.1.

Finally, the overall assessment score \mathcal{S}_t and its corresponding textual explanation \mathcal{T}_t are derived by aggregating the evaluations for query novelty and usefulness. The overall score is determined by a logical conjunction of the component scores:

$$\mathcal{S}_t = \begin{cases} 1, & \text{if } \mathcal{S}_t^{\text{novel}} = 1 \wedge \mathcal{S}_t^{\text{useful}} = 1, \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

The final explanation is synthesized by concatenating the textual feedback from both components:

$$\mathcal{T}_t = \mathcal{T}_t^{\text{novel}} \parallel \mathcal{T}_t^{\text{useful}}, \quad (12)$$

where \parallel denotes the concatenation operator.

4.3 Process Reward-Guided Query Refinement

This section introduces the query refinement mechanism, which is designed to promote the optimization of query generation. It is achieved by systematically identifying and refining low-quality queries, then regenerating subsequent search steps from these refined points. This mechanism serves a pivotal function within the three-stage curriculum learning framework by generating comparative data for training and acting as a rollout strategy.

Formally, this process can be represented as follows. The search agent starts by generating a complete trajectory H_T , represented as $(q, \tau_0, a_0, o_0, \dots, \tau_i, a_i, o_i, \dots, \tau_T, a_T)$. Each search query in this trajectory is then evaluated by the process rewards mechanism,

yielding a sequence of scores $(\mathcal{S}_0, \mathcal{S}_1, \dots, \mathcal{S}_{T-1})$ and corresponding textual explanations $(\mathcal{T}_0, \mathcal{T}_1, \dots, \mathcal{T}_{T-1})$ with Equation (11) and (12). For each low-quality query a_i where the score $\mathcal{S}_i = 0$, a refinement step is triggered. The refined query a'_i is generated by a language model as follows:

$$a'_i = \text{LLM}_{\text{refine}}(q, H_i, \mathcal{T}_i). \quad (13)$$

Here, $\text{LLM}_{\text{refine}}$ is the same lightweight SFT-tuned model introduced earlier, q is the user's original query, H_i is the trajectory history up to step i , and \mathcal{T}_i is the textual feedback diagnosing the quality issue for the low-quality query a_i . The search agent subsequently regenerates the search process from this refined query a'_i , yielding a new trajectory H'_T , represented as $(q, \tau_0, a_0, o_0, \dots, \tau_i, a'_i, o'_i, \dots, \tau'_T, a'_T)$. The primary distinction between the initial and revised trajectories originates from the refined query a'_i , resulting in a different reward for a_i and a'_i , thereby promoting the optimization of query generation within the curriculum learning framework.

Specifically, to enable the model to effectively refine the low-quality search query based on the textual feedback provided by the process rewards mechanism, we distill key empirical insights into a set of actionable guidelines based on a thorough analysis of representative cases:

- If the textual feedback indicates that the query is redundant or unnecessary, the refined query should serve for a more necessary intent and eliminate redundancy.
- If the textual feedback indicates that the retrieved results do not contain the expected information, the model should strategically reformulate the query to better capture the target content. This

reformulation may involve switching between a complete semantic question and a keyphrase-based query, or adaptively adding or removing information from the original query.

4.4 Query-Oriented Training Framework

This section presents a three-stage curriculum learning framework that integrates the two preceding mechanisms, enabling the agent to progressively internalize the ability to improve query quality through a progression from imitation, to alignment, and ultimately to generalization. The following paragraphs detail the three progressive stages: Query Quality Screened Imitation Learning, Query Generation Alignment, and Query Aware Policy Optimization.

Stage-1: Query Quality Screened Imitation Learning. In this stage, we employ SFT to guide the model in its initial learning of information retrieval and utilization. A critical step in SFT is the selection of high-quality trajectories for training. Following common practice, we begin by selecting trajectories that yield correct final answers and adhere to the proper format, thus guiding the model towards correct patterns from the outset. However, many trajectories, despite yielding correct final answers, contain low-quality intermediate search queries. Learning from such trajectories could lead the model to pick up suboptimal behaviors, thereby impairing its overall performance. To address this, we further leverage process rewards to selectively retain only those trajectories comprised entirely of high-quality intermediate search queries, i.e., $\forall t \in [0, \dots, T], S_t = 1$. This ensures that the trajectories comprising our final training dataset \mathcal{D} not only yield correct final answers but also exhibit high-quality intermediate search queries. We then apply the standard SFT objective, which is formulated as:

$$\mathcal{L}_{\text{SFT}}(\theta) = -\mathbb{E}_{(q,y) \sim \mathcal{D}} [\log P_\theta(y | q)], \quad (14)$$

where q is the user’s original query, y is the agent’s high-quality response, and θ denotes the model parameters.

Stage-2: Query Generation Alignment. In this stage, the search agent cultivates advanced query generation capabilities through DPO training. Unlike common approaches that directly generate trajectories from scratch, we employ the query refinement mechanism when constructing comparative data. For each user’s query q , the search agent first generates an initial trajectory y_0 . Following this, each low-quality query within y_0 is refined and the search agent regenerates subsequent search steps from the refined query, producing a sequence of trajectories y_1, \dots, y_n , where n is the number of low-quality queries. This process ensures that for a given input q , the key differences among the candidate trajectories y_0, y_1, \dots, y_n originate specifically from the refined queries, thereby directly promoting the optimization of query generation.

Next, for each user’s query q , we choose one positive sample y_w and one negative sample y_l among the corresponding candidate trajectories y_0, y_1, \dots, y_n . Diverging from approaches that rely only on the correctness of the final answer, our selection criteria incorporate *both the final-answer correctness and the quality of intermediate search queries*, guided by the following principles:

- A trajectory with a correct final answer is preferred over one with an incorrect answer.

- Among trajectories with correct final answers, those with fewer low-quality (i.e., $S_t = 0$) queries are preferred.
- Among trajectories with incorrect final answers, those containing more high-quality (i.e., $S_t = 1$) queries are preferred.

We then optimize the model using the standard DPO objective:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(q,y_w,y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_w | q)}{\pi_{\text{ref}}(y_w | q)} - \beta \log \frac{\pi_\theta(y_l | q)}{\pi_{\text{ref}}(y_l | q)} \right) \right]. \quad (15)$$

Here, q is the user’s original query, y_w is the positive sample, y_l is the negative sample, β represents the hyperparameter, σ refers to the sigmoid function, θ is the model parameters, and π_{ref} indicates the reference model, which is initialized to π_θ and kept frozen during training.

Stage-3: Query Aware Policy Optimization. In the final stage, we further enhance the search agent’s integrated capabilities of information retrieval and utilization through Query Aware Policy Optimization. Specifically, we train it on a curated set of challenging questions that remained unresolved after multiple sampling trials. Unlike the standard GRPO algorithm that generates G independent trajectories from scratch, our method employs the query refinement mechanism as its rollout strategy. For each user’s query, the search agent first generates an initial trajectory y_0 and then expands it into y_0, y_1, \dots, y_n through sequential refinement and regeneration. Different from the Query Generation Alignment stage, we retain at most M trajectories from this set to avoid too many trajectories sharing a common prefix, thereby ensuring behavioral diversity and promoting the holistic improvement of the agent’s capabilities. If the total number of trajectories collected remains less than G , we repeat this generation-and-expansion process, until a complete set of G trajectories is obtained.

For reward design, we integrate process supervision into the reward function. Following Eq. (6), our reward function is:

$$r = r_{\text{composite}} + \lambda \cdot r_{\text{format}}, \quad (16)$$

where λ is a weighting coefficient, $r_{\text{format}} \in \{0, 1\}$ indicates the correctness of the output format, and $r_{\text{composite}}$ defined in Eq. (7) integrates both outcome and process reward as follows:

$$r_{\text{composite}} = \begin{cases} \max(r_{\text{outcome}} - \gamma \cdot n_{\text{wrong}}, \phi_{\min}), & r_{\text{outcome}} = 1, \\ \min(r_{\text{outcome}} + \gamma \cdot n_{\text{correct}}, \phi_{\max}), & r_{\text{outcome}} = 0. \end{cases} \quad (17)$$

Here, $r_{\text{outcome}} \in \{0, 1\}$ denotes the final answer’s correctness, n_{wrong} and n_{correct} represent the number of low- (i.e., $S_t = 0$) and high-quality (i.e., $S_t = 1$) queries respectively, γ is a scaling factor for process rewards, and ϕ_{\min}, ϕ_{\max} bound the influence of process rewards. This reward design incentivizes the agent not only to prioritize final answer correctness but also to refine its search process by reducing low-quality queries in successful trajectories. Moreover, even when unable to provide a final correct answer, the agent is motivated to generate more high-quality queries that may progressively approach the solution. We then optimize the model using the standard GRPO objective introduced in Section 3.2.

Table 1: Performance comparison of SmartSearch and existing approaches on four knowledge-intensive benchmarks, with bold for the best and underlined for the runner-up. Numbers in () indicate the improvement compared with the runner-up.

Method	2WikiMQA		HotpotQA		Bamboogle		Musique		Average	
	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1
<i>Prompt-based Approaches</i>										
Direct Inference	19.3	24.7	14.6	24.5	4.0	11.6	2.3	7.9	10.1	17.2
CoT	18.1	24.9	12.8	24.0	14.4	25.5	2.2	7.8	11.9	20.6
IRCoT	20.0	27.2	19.3	28.0	16.8	25.9	5.8	12.7	15.5	23.5
RAG	22.5	31.4	24.3	36.7	7.2	17.2	4.5	12.2	14.6	24.4
Search-o1	20.9	29.4	22.0	33.6	28.8	36.1	5.1	12.6	19.2	27.9
<i>RL Approaches with Outcome Rewards</i>										
ReSearch	29.4	36.7	28.5	40.8	12.8	22.9	10.0	17.3	20.2	29.4
ZeroSearch	29.2	36.5	27.5	39.1	14.4	25.4	10.4	18.2	20.4	29.8
R1-Searcher	29.8	37.1	27.0	38.7	31.2	39.2	8.0	16.4	24.0	32.9
Search-R1	27.3	35.5	31.9	41.1	29.4	38.8	9.3	16.6	24.5	33.0
<i>RL Approaches with Process Rewards</i>										
ReasonRag	<u>36.5</u>	<u>43.2</u>	32.2	41.7	30.4	39.1	11.3	18.6	27.6	35.7
PPR	33.7	41.8	<u>38.1</u>	<u>50.3</u>	31.2	39.4	14.7	22.0	29.4	38.4
StepSearch	32.1	38.9	35.1	45.9	<u>36.8</u>	<u>48.4</u>	<u>16.6</u>	<u>24.9</u>	<u>30.1</u>	<u>39.5</u>
SmartSearch (Ours)	45.3(↑24%)	52.3(↑21%)	40.7(↑7%)	52.4(↑4%)	44.8(↑22%)	56.1(↑16%)	19.1(↑15%)	27.8(↑12%)	37.5(↑25%)	47.2(↑19%)

5 Experiments

5.1 Experimental Setup

Dataset. We comprehensively assess SmartSearch’s performance through experiments on two types of benchmarks: (1) knowledge-intensive tasks, including 2WikiMultihopQA [16], HotpotQA [61], Bamboogle [33], and Musique [47], and (2) web exploration tasks, including GAIA [32] and WebWalker [55].

Metrics. For a consistent comparison with previous studies, we use the widely adopted Exact Match (EM) and word-level F1 score to assess the answers’ correctness. To assess search efficiency, we follow prior work [5] and employ the Search Efficiency metric, defined as: $S_E = \frac{1}{N} \sum_{i=1}^N \frac{F_i}{T_i}$. Here, N represents the dataset size, F_i denotes the F1 score for sample i , and T_i represents the search call count for sample i . Additionally, to assess search query quality, we introduce the Search Quality metric, defined as: $S_Q = \frac{1}{N} (C_{\text{perfect}} + C_{\text{partial}})$ where N represents the dataset size, C_{perfect} denotes the number of samples where the final answer is correct and all intermediate search queries are of high quality, and C_{partial} denotes the number of samples where the final answer is incorrect but the trajectory contains high-quality intermediate search queries. In particular, we define the Perfect Rate as $\frac{1}{N} C_{\text{perfect}}$ and the Partial Rate as $\frac{1}{N} C_{\text{partial}}$, which contribute to the overall Search Quality metric from two different aspects.

Baselines. We compare SmartSearch with several representative baselines, which are classified into three categories: (1) prompt-based approaches, including Direct Inference, CoT [53], IRCoT [48], RAG [27], and Search-o1 [28]. (2) RL approaches with outcome rewards, including ReSearch [4], ZeroSearch [40], R1-Searcher [39], and Search-R1 [21]. (3) RL approaches with process rewards, including PPR [59], ReasonRag [67], and StepSearch [52].

Table 2: Performance comparison of SmartSearch and baselines on web exploration tasks, with bold for the best.

Method	GAIA		WebWalker		Average	
	EM	F1	EM	F1	EM	F1
Search-o1	4.7	8.0	6.3	18.3	5.5	13.2
Search-R1	6.3	9.8	7.5	21.6	6.9	15.7
StepSearch	9.4	12.5	9.1	25.8	9.3	19.2
SmartSearch	13.4	16.7	11.5	31.0	12.5	23.9

Implementation Details. Qwen2.5-3B-Instruct serves as the base model in SmartSearch and other baselines. For local search, we utilize the 2018 Wikipedia dump [24] provided by FlashRAG [22] as the knowledge corpus, and employ E5-base-v2 [50] as the retriever to obtain top 5 documents. For web search, the Serper API is employed to retrieve the top 10 document snippets. Our training is conducted under the LLaMA-Factory and VERL frameworks, using Asearcher-Base as the training dataset. Additionally, to improve efficiency, we train a smaller student model, Qwen2.5-3B-Instruct, to perform scoring and query refinement, with training labels annotated by the teacher model, Qwen3-32B.

5.2 Main Result

Overall Performance. Table 1 presents the main results, demonstrating that SmartSearch consistently surpasses existing approaches across four datasets, and yielding several important insights.

(1) Prompt-based approaches exhibit limited performance. Direct Inference and CoT, which are based entirely on the model’s internal knowledge, achieve an average EM of only around 10%,

Table 3: Ablation study results for the two core mechanisms in SmartSearch across all three stages of curriculum learning training framework, with bold for the best results of each stage.

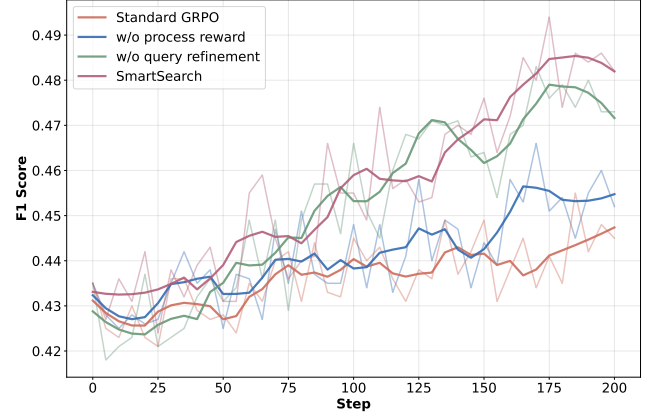
Method	2WikiMQA		HotpotQA		Bamboogle		Musique		Average	
	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1
Stage 1										
SmartSearch	38.2	45.3	35.3	45.5	38.4	51.0	14.7	21.6	31.7	40.9
w/o process rewards	33.8	40.6	32.5	42.8	36.0	48.7	12.6	20.8	28.7	38.2
Stage 2										
SmartSearch	41.4	48.7	37.9	48.5	39.2	51.8	15.4	23.6	33.5	43.2
w/o process rewards	40.2	47.4	36.5	47.2	37.6	50.1	14.4	22.3	32.2	41.8
w/o query refinement	39.2	46.7	35.6	46.1	36.0	49.6	14.6	22.9	31.4	41.3
Stage 3										
SmartSearch	45.3	52.3	40.7	52.4	44.8	56.1	19.1	27.8	37.5	47.2
w/o process rewards	43.3	50.6	40.0	51.5	39.2	51.6	17.9	26.7	35.1	45.1
w/o query refinement	44.1	51.2	39.9	51.6	41.6	54.2	17.5	26.4	35.8	45.9
Standard GRPO	43.6	50.8	39.0	50.6	40.8	52.5	15.9	24.8	34.8	44.7

highlighting the inherent challenges of LLMs, including hallucinations and static parametric knowledge. RAG and IRCot yield a gain of around 5% in average EM, demonstrating the necessity of integrating external knowledge. Among prompt-based methods, Search-o1 attains the highest performance, reaching an average EM score of 19.2%, reflecting the effectiveness of search agents. However, it still lags behind other fine-tuning-based approaches.

(2) Incorporating process rewards effectively enhances RL training. While outcome-driven RL methods such as ReSearch, Search-R1, and R1-Searcher improve performance over prompt-based approaches, indicating the benefits of RL for LLM-based search agents, they often remain inferior to RL approaches that integrate both outcome and process rewards by a margin of around 5% in both average EM and F1 score. This underscores that reward signals based solely on final outcomes result in sparse feedback. Such sparse feedback provides insufficient guidance for intermediate steps and leads to unstable optimization, thereby highlighting the critical importance of fine-grained supervision.

(3) Optimizing the quality of intermediate search queries significantly improves overall performance. Existing methods often overlook the quality of intermediate search queries, which can lead to stagnation in information retrieval abilities. By explicitly optimizing the quality of intermediate search queries under the guidance of process rewards, SmartSearch enhances the search agent’s overall effectiveness, achieving more than 7% improvement in both average EM and F1 score compared with other process-supervised RL methods.

Generalization to Web Search Scenarios. As discussed earlier, SmartSearch is trained solely on Wikipedia-based local search. To evaluate its generalization ability to web search, we test it against several baseline models on two demanding web exploration tasks, GAIA and WebWalker. As demonstrated in Table 2, SmartSearch surpasses existing approaches across both datasets, achieving an average F1 score increase of nearly 5%. This indicates that while SmartSearch optimizes the quality of intermediate search queries

**Figure 4: F1 score training dynamics for different algorithms.**

in the local search setting, it also exhibits strong generalization capabilities in the web search environment, maintaining robust performance despite the difference in settings.

5.3 Ablation Study

To further examine the impact of SmartSearch’s two key mechanisms—process rewards and query refinement, we conduct extensive ablation studies across all three training stages. The results are summarized in Table 3.

For Stage 1, we compare our configuration with a baseline that filters the training data exclusively according to whether the final answer is correct. The results indicate that incorporating query-quality filtering enables the model to achieve superior performance with only 60% of the training data, highlighting the importance of learning from trajectories with high-quality search processes.

For Stage 2, we compare our method with two alternatives: (1) directly generating full trajectories without query refinement, and (2) determining preference based exclusively on the final answer

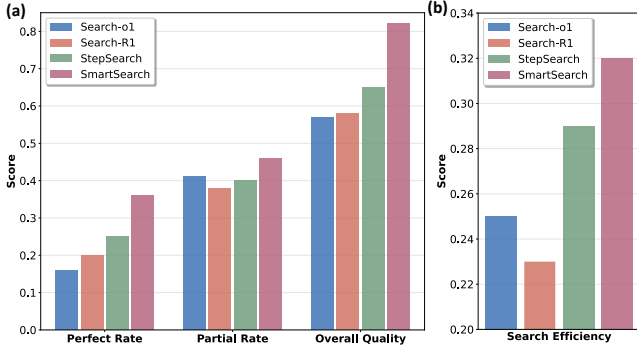


Figure 5: Left: Search query quality comparison. Right: Search efficiency comparison.

correctness. Ablation results underscore the essential contribution of both mechanisms in this stage, particularly the query refinement mechanism, underscoring the significance of promoting the optimization of query generation.

For Stage 3, we compare our algorithm with three variants: a GRPO baseline, a version that only incorporates the process rewards into the reward function, and a version that only applies the query refinement mechanism during rollout. As depicted in Figure 4, we demonstrate F1 score curves of various RL algorithms during training. The results demonstrate that our algorithm reliably outperforms all alternatives. Notably, integrating process rewards into the reward function yields significant gains, illustrating the crucial role of fine-grained supervision for the quality of each query.

5.4 Quantitative Analyses

To comprehensively assess the effectiveness of the SmartSearch framework, we perform multiple quantitative experiments. The following analyses demonstrate its superiority in four key aspects: intermediate query quality, search efficiency, the effectiveness of the process reward model, and the effectiveness-efficiency trade-off.

Search Query Quality Analysis. To assess whether SmartSearch improves the quality of intermediate search queries, we compare the Search Quality metric across various methods. As presented in Figure 5 (a), SmartSearch achieves the highest Search Quality, with the highest values for both Perfect Rate and Partial Rate, which contribute to the overall Search Quality metric. This indicates that SmartSearch effectively enhances the quality of intermediate search queries. Specifically, the search agent reduces ineffective searches while striving to generate perfect trajectories where all queries are of high quality. Even when unable to provide a final correct answer, the agent makes more attempts to generate high-quality queries that edge closer to the correct solution.

Search Efficiency Analysis. The results in previous sections have shown that SmartSearch outperforms all baselines in accuracy. We now further evaluate whether it also achieves superior search efficiency. To this end, we compare the search efficiency metrics across multiple methods, and as shown in Figure 5 (b), SmartSearch outperforms all other methods in this regard. This suggests that by optimizing the quality of intermediate search queries, SmartSearch

generates more precise queries, reducing ineffective or failed search rounds and, as a result, improving overall search efficiency.

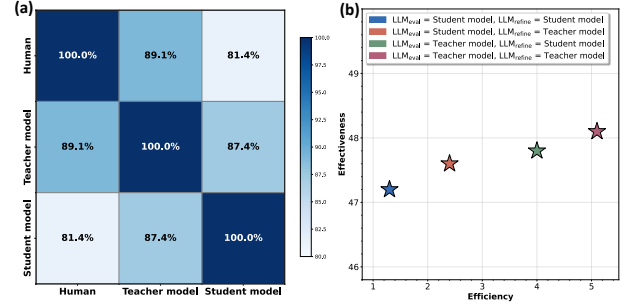


Figure 6: Left: Overlap between the scores assigned to queries by the student model, teacher model, and human annotations. Right: Effectiveness and efficiency tradeoff in SmartSearch.

Process Reward Model Analysis. The process reward model plays a crucial role in our approach by providing fine-grained supervision for the quality of each query and guiding subsequent query refinement. To assess the effectiveness of our process reward model, we randomly choose 100 trajectories. For each search query in these trajectories, we compare the scores annotated by the teacher model, the student model, and human annotators. Figure 6 (a) illustrates the overlap between the scores assigned to each intermediate search query by these three sources. The results reveal that the teacher model achieves nearly 90% overlap with human annotations, demonstrating its effectiveness in labeling the training data. After training, the student model achieves over 85% overlap with the teacher model, indicating the effectiveness of our fine-tuning. Finally, the student model shows over 80% overlap with human annotations, a result that is entirely acceptable, striking a good balance between scoring accuracy and efficiency.

Effectiveness-Efficiency Trade-off. To validate the suitability of the lightweight student model as both LLM_{eval} and LLM_{refine}, we conduct experiments using a more powerful teacher model in these roles. In this experiment, effectiveness is defined as the average F1 score, while efficiency refers to the average time (s) required to apply the process rewards and query refinement mechanisms to each sample. As shown in Figure 6 (b), using a more powerful model as both LLM_{eval} and LLM_{refine} indeed improves the agent’s effectiveness, underscoring the importance of the process rewards and query refinement mechanisms within our training framework. However, this gain in average F1 score is modest, remaining below 1%, whereas the time required to process each sample increases by nearly five times. This result demonstrates a clear trade-off between effectiveness and efficiency, and the decision to use the lightweight student model as both LLM_{eval} and LLM_{refine} proves to be a sensible one. This choice achieves an optimal balance between effectiveness and efficiency, ensuring effective query optimization while avoiding excessive computational costs.

6 Conclusion

In this work, we introduce SmartSearch, a framework designed to optimize the quality of intermediate search queries through two key mechanisms: (1) Process rewards, which provide fine-grained supervision for the quality of each query through Dual-Level Assessment. (2) Query refinement, which promotes the optimization of query generation by selectively refining low-quality queries and regenerating subsequent search rounds from these refined points. Building on the foundation of the two mechanisms, we design a three-stage curriculum learning framework that guides the agent through a progression from imitation and alignment to generalization, enabling it to progressively internalize the ability to enhance query quality. Experiments across four challenging benchmarks demonstrate that SmartSearch consistently surpasses existing baselines, with further quantitative analyses confirming significant gain in both search efficiency and query quality.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [3] Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 17754–17762.
- [4] Mingyang Chen, Linzhuang Sun, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Haofer Wang, Jeff Z Pan, Wen Zhang, Huajun Chen, et al. 2025. Learning to reason with search for llms via reinforcement learning. *arXiv preprint arXiv:2503.19470* (2025).
- [5] Yifei Chen, Guanting Dong, and Zhicheng Dou. 2025. Toward Effective Tool-Integrated Reasoning via Self-Evolved Preference Learning. *arXiv preprint arXiv:2509.23285* (2025).
- [6] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research* 24, 240 (2023), 1–113.
- [7] Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. [n. d.]. SFT Memorizes, RL Generalizes: A Comparative Study of Foundation Model Post-training. In *Forty-second International Conference on Machine Learning*.
- [8] Tri Dao. 2024. FLASHATTENTION-2: FASTER ATTENTION WITH BETTER PARALLELISM AND WORK PARTITIONING. In *12th International Conference on Learning Representations, ICLR 2024*.
- [9] Yong Deng, Guoqing Wang, Zhenzhe Ying, Xiaofeng Wu, Jinzhen Lin, Wenwen Xiong, Yuqin Dai, Shuo Yang, Zhanwei Zhang, Qiwen Wang, et al. 2025. Atom-searcher: Enhancing agentic deep research via fine-grained atomic thought reward. *arXiv preprint arXiv:2508.12800* (2025).
- [10] Guanting Dong, Yifei Chen, Xiaoxi Li, Jiajie Jin, Hongjin Qian, Yutao Zhu, Hangyu Mao, Guorui Zhou, Zhicheng Dou, and Ji-Rong Wen. 2025. Tool-Star: Empowering LLM-Brained Multi-Tool Reasoner via Reinforcement Learning. *arXiv preprint arXiv:2505.16410* (2025).
- [11] Guanting Dong, Hangyu Mao, Kai Ma, Licheng Bao, Yifei Chen, Zhongyuan Wang, Zhongxia Chen, Jiazhen Du, Huiyang Wang, Fuzheng Zhang, et al. 2025. Agentic reinforced policy optimization. *arXiv preprint arXiv:2507.19849* (2025).
- [12] Guanting Dong, Yutao Zhu, Chenghao Zhang, Zechen Wang, Ji-Rong Wen, and Zhicheng Dou. 2025. Understand what LLM needs: Dual preference alignment for retrieval-augmented generation. In *Proceedings of the ACM on Web Conference 2025*. 4206–4225.
- [13] Tianqing Fang, Zhisong Zhang, Xiaoyang Wang, Rui Wang, Can Qin, Yuxuan Wan, Jun-Yu Ma, Ce Zhang, Jiaqi Chen, Xiyun Li, et al. 2025. Cognitive kernel-pro: A framework for deep research agents and agent foundation models training. *arXiv preprint arXiv:2508.00414* (2025).
- [14] Jiaxuan Gao, Wei Fu, Minyang Xie, Shusheng Xu, Chuyi He, Zhiyu Mei, Banghua Zhu, and Yi Wu. [n. d.]. Beyond Ten Turns: Unlocking Long-Horizon Agentic Search with Large-Scale Asynchronous RL. In *First Workshop on Multi-Turn Interactions in Large Language Models*.
- [15] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofer Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997* (2023).
- [16] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing A Multi-hop QA Dataset for Comprehensive Evaluation of Reasoning Steps. In *Proceedings of the 28th International Conference on Computational Linguistics*. 6609–6625.
- [17] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *Comput. Surveys* 55, 12 (2023), 1–38.
- [18] Jinhao Jiang, Jiayi Chen, Junyi Li, Ruiyang Ren, Shijie Wang, Wayne Xin Zhao, Yang Song, and Tao Zhang. 2025. Rag-star: Enhancing deliberative reasoning with retrieval augmented verification and refinement. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. 7064–7074.
- [19] Pengcheng Jiang, Jiacheng Lin, Lang Cao, Runchu Tian, SeongKu Kang, Zifeng Wang, Jimeng Sun, and Jiawei Han. 2025. DeepRetrieval: Hacking Real Search Engines and Retrievers with Large Language Models via Reinforcement Learning. *CoRR* (2025).
- [20] Yi Jiang, Lei Shen, Lujie Niu, Sendong Zhao, Wenbo Su, and Bo Zheng. 2025. QAgent: A modular Search Agent with Interactive Query Understanding. *arXiv preprint arXiv:2510.08383* (2025).
- [21] Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. [n. d.]. Search-R1: Training LLMs to Reason and Leverage Search Engines with Reinforcement Learning. ([n. d.]).
- [22] Jiajie Jin, Yutao Zhu, Zhicheng Dou, Guanting Dong, Xinyu Yang, Chenghao Zhang, Tong Zhao, Zhao Yang, and Ji-Rong Wen. 2025. Flashrag: A modular toolkit for efficient retrieval-augmented generation research. In *Companion Proceedings of the ACM on Web Conference 2025*. 737–740.
- [23] Ehsan Kamaloo, Nouha Dziri, Charles Clarke, and Davood Rafiei. 2023. Evaluating Open-Domain Question Answering in the Era of Large Language Models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 5591–5606.
- [24] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick SH Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *EMNLP (1)*. 6769–6781.
- [25] Shanglin Lei, Guanting Dong, Xiaoping Wang, Keheng Wang, and Sirui Wang. 2023. Instructor: Reforming emotion recognition in conversation with a retrieval multi-task llms framework. *CoRR* (2023).
- [26] Yongqi Leng, Yikun Lei, Xikai Liu, Meizhi Zhong, Bojian Xiong, Yurong Zhang, Yan Gao, Yao Hu, Deyi Xiong, et al. 2025. DecEx-RAG: Boosting Agentic Retrieval-Augmented Generation with Decision and Execution Optimization via Process Supervision. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: Industry Track*. 1412–1425.
- [27] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems* 33 (2020), 9459–9474.
- [28] Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. 2025. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366* (2025).
- [29] Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yongkang Wu, Ji-Rong Wen, Yutao Zhu, and Zhicheng Dou. 2025. Webthinker: Empowering large reasoning models with deep research capability. *arXiv preprint arXiv:2504.21776* (2025).
- [30] Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yongkang Wu, Zhonghua Li, Ye Qi, and Zhicheng Dou. 2025. Retrollm: Empowering large language models to retrieve fine-grained evidence within generation. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 16754–16779.
- [31] Pan Lu, Bowen Chen, Sheng Liu, Rahul Thapa, Joseph Boen, and James Zou. [n. d.]. OctoTools: An Agentic Framework with Extensible Tools for Complex Reasoning. In *ICLR 2025 Workshop on Foundation Models in the Wild*.
- [32] Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. [n. d.]. Gaia: a benchmark for general ai assistants. In *The Twelfth International Conference on Learning Representations*.
- [33] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. 5687–5711.
- [34] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. DeepSpeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*. 3505–3506.
- [35] Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922* (2023).
- [36] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *Advances in*

- Neural Information Processing Systems* 36 (2023), 68539–68551.
- [37] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300* (2024).
- [38] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. 2023. Towards expert-level medical question answering with large language models. *arXiv preprint arXiv:2305.09617* (2023).
- [39] Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. 2025. R1-Searcher: Incentivizing the Search Capability in LLMs via Reinforcement Learning. *CoRR* (2025).
- [40] Hao Sun, Zile Qiao, Jiayan Guo, Xuanbo Fan, Yingyan Hou, Yong Jiang, Pengjun Xie, Yan Zhang, Fei Huang, and Jingren Zhou. 2025. Zerosearch: Incentivize the search capability of llms without searching. *arXiv preprint arXiv:2505.04588* (2025).
- [41] Liyan Tang, Zhaoyi Sun, Betina Idnay, Jordan G Nestor, Ali Soroush, Pierre A Elias, Ziyang Xu, Ying Ding, Greg Durrett, Justin F Rousseau, et al. 2023. Evaluating large language models on medical evidence summarization. *NPJ digital medicine* 6, 1 (2023), 158.
- [42] Zhengwei Tao, Haiyang Shen, Baixuan Li, Wenbiao Yin, Jialong Wu, Kuan Li, Zhongwang Zhang, Huifeng Yin, Rui Ye, Liwen Zhang, et al. 2025. Webleaper: Empowering efficiency and efficacy in webagent via enabling info-rich seeking. *arXiv preprint arXiv:2510.24697* (2025).
- [43] Kimi Team. [n. d.]. KIMI K2: OPEN AGENTIC INTELLIGENCE. ([n. d.]).
- [44] Qwen Team. 2024. Qwq: Reflect deeply on the boundaries of the unknown. *Hugging Face* (2024).
- [45] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
- [46] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).
- [47] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. MuSiQue: Multihop Questions via Single-hop Question Composition. *Transactions of the Association for Computational Linguistics* 10 (2022), 539–554.
- [48] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. In *Proceedings of the 61st annual meeting of the association for computational linguistics (volume 1: long papers)*. 10014–10037.
- [49] Guoqing Wang, Sunhao Dai, Guangze Ye, Zeyu Gan, Wei Yao, Yong Deng, Xiaofeng Wu, and Zhenzhe Ying. 2025. Information Gain-based Policy Optimization: A Simple and Effective Approach for Multi-Turn LLM Agents. *arXiv preprint arXiv:2510.14967* (2025).
- [50] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533* (2022).
- [51] Yiding Wang, Zhepei Wei, Xinyu Zhu, and Yu Meng. 2025. Beyond Outcome Reward: Decoupling Search and Answering Improves LLM Agents. *arXiv preprint arXiv:2510.04695* (2025).
- [52] Ziliang Wang, Xuhui Zheng, Kang An, Cijun Ouyang, Jialu Cai, Yuhang Wang, and Yichao Wu. 2025. StepSearch: Igniting LLMs Search Ability via Step-Wise Proximal Policy Optimization. *arXiv preprint arXiv:2505.15107* (2025).
- [53] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [54] Tongyu Wen, Chenglong Wang, Xiyuan Yang, Haoyu Tang, Yueqi Xie, Lingjuan Lyu, Zhicheng Dou, and Fangzhao Wu. 2025. Defending against Indirect Prompt Injection by Instruction Detection. *CoRR abs/2505.06311* (2025). arXiv:2505.06311 doi:10.48550/ARXIV.2505.06311
- [55] Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He, Deyu Zhou, Pengjun Xie, et al. [n. d.]. WebWalker: Benchmarking LLMs in Web Traversal. In *Workshop on Reasoning and Planning for Large Language Models*.
- [56] Peilin Wu, Mian Zhang, Kun Wan, Wentian Zhao, Kaiyu He, Xinya Du, and Zhiyu Chen. 2025. HiPRAG: Hierarchical Process Rewards for Efficient Agentic Retrieval Augmented Generation. *arXiv preprint arXiv:2510.07794* (2025).
- [57] Guangzhi Xiong, Qiao Jin, Xiao Wang, Yin Fang, Haolin Liu, Yifan Yang, Fangyuan Chen, Zhixing Song, Dengyu Wang, Minjia Zhang, et al. 2025. Rag-gym: Optimizing reasoning and search agents with process supervision. *arXiv preprint arXiv:2502.13957* (2025).
- [58] Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. 2023. A paradigm shift in machine translation: Boosting translation performance of large language models. *arXiv preprint arXiv:2309.11674* (2023).
- [59] Peiran Xu, Zhuohao Li, Xiaoying Xing, Guannan Zhang, Debiao Li, and Kunyu Shi. 2025. Hybrid Reward Normalization for Process-supervised Non-verifiable Agentic Tasks. *arXiv preprint arXiv:2509.25598* (2025).
- [60] Zhichao Xu, Zongyu Wu, Yun Zhou, Aosong Feng, Kang Zhou, Sangmin Woo, Kiran Ramnath, Yijun Tian, Xuan Qi, Weikang Qiu, et al. 2025. Beyond Correctness: Rewarding Faithful Reasoning in Retrieval-Augmented Generation. *arXiv preprint arXiv:2510.13272* (2025).
- [61] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 conference on empirical methods in natural language processing*. 2369–2380.
- [62] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*.
- [63] Siliang Zeng, Quan Wei, William Brown, Oana Frunza, Yuriy Nevmyvaka, Yang Katie Zhao, and Mingyi Hong. 2025. Reinforcing multi-turn reasoning in llm agents via turn-level credit assignment. In *ICML 2025 Workshop on Computer Use Agents*.
- [64] Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study. In *International Conference on Machine Learning*. PMLR, 41092–41110.
- [65] Guibin Zhang, Hejia Geng, Xiaohang Yu, Zhenfei Yin, Zaibin Zhang, Zelin Tan, Heng Zhou, Zhongzhi Li, Xiangyuan Xue, Yijiang Li, et al. 2025. The landscape of agentic reinforcement learning for llms: A survey. *arXiv preprint arXiv:2509.02547* (2025).
- [66] Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori B Hashimoto. 2024. Benchmarking large language models for news summarization. *Transactions of the Association for Computational Linguistics* 12 (2024), 39–57.
- [67] Wenlin Zhang, Xiangyang Li, Kuicai Dong, Yichao Wang, Pengyue Jia, Xiaopeng Li, Yingyi Zhang, Derong Xu, Zhaocheng Du, Huifeng Guo, et al. 2025. Process vs. Outcome Reward: Which is Better for Agentic RAG Reinforcement Learning. *arXiv preprint arXiv:2505.14069* (2025).
- [68] Qingfei Zhao, Ruobing Wang, Dingling Xu, Daren Zha, and Limin Liu. 2025. R-Search: Empowering LLM Reasoning with Search via Multi-Reward Reinforcement Learning. *arXiv preprint arXiv:2506.04185* (2025).

A Prompt Templates

Prompt for SmartSearch.

You are a helpful assistant that can solve the given question step by step with the help of the wikipedia search tool. Given a question, you need to first think about the reasoning process in the mind and then provide the answer. During thinking, you can invoke the wikipedia search tool to search for fact information about specific topics if needed. The reasoning process and answer are enclosed within <think> </think> and <answer> </answer> tags respectively, and the search query and result are enclosed within <search> </search> and <result> </result> tags respectively. For example, <think> This is the reasoning process. </think> <search> search query here </search> <result> search result here </result> <think> This is the reasoning process. </think> <answer> The final answer is

answer here

</answer>. In the last part of the answer, the final exact answer is enclosed within \boxed{ } with latex format.

Prompt for scoring.

You are a query-evaluation assistant. Your task is to assess the quality of a search agent’s query of the current search round according to the user’s question, the golden answer and the agent’s search process up to the current search round.

If the agent’s query intent of the current search round is necessary, and the corresponding query result includes the answer for the query, the score for query should be 1. Otherwise, the score for the query should be 0. The details of the assessment are in the Evaluation Guideline, please read it carefully.

User’s question

{question}

Golden answer

{answer}

Agent’s search process up to the current search round

{context}

Evaluation Guideline

1. Identify the agent’s query intent of the current search round accurately.
2. The query result doesn’t need to solve the user’s question directly; but it must include the information that address the agent’s query intent completely, related entities alone is not enough.
3. The intended entity and the one found in the query result must be exactly the same, otherwise, the score should be 0.

Output Format:

<answer> score for the query </answer>

<explanation> explanation for the score </explanation>

Prompt for Refining.

You are a query-refine assistant. Your task is to refine a search agent’s query of the current search round within <search> </search> according to the user’s question, the agent’s search process up to the current search round and the issues of the query. The details of the refinement are in the Refine Guideline, please read it carefully.

User’s question

{question}

Agent’s search process up to the current search round

{context}

Issues of the query

{explanation}

Refine Guideline

1. The refined query is meant to replace the query of the current round, so don’t rely on any query result within <result> </result> from the current round when refining the query.
2. If the issues of the query indicate that the query intent is unreasonable, the refined query should serve for a more necessary and actionable query intent.
3. The refined query can be expressed as a complete semantic question or a keyphrase-based query, and you may add or remove information from the original query. All depends on which option best serves the agent’s query intent, ensuring that the query result contains the answer to the agent’s query intent.

Output format:

<search> refined query </search>

<explanation> explanation for the refined query </explanation>

B Implementation Details

In the Query Quality Screened Imitation Learning stage, we employ ARPO-14B [11] as the policy model for trajectory sampling. The trajectories obtained through this sampling process are then used to fine-tune the Qwen2.5-3B-Instruct model through SFT, resulting in the SFT model. The training is conducted with a learning rate of $7e-6$ over a total of 3 epochs, and the maximum input length during training is set to 16384 tokens. We utilize DeepSpeed ZeRO-3 [34] and FlashAttention2 [8] to accelerate training, with a total batch size of 64 and applying BF16 precision.

In the Query Generation Alignment stage, we perform DPO training using trajectories generated by the SFT model as positive and negative samples, resulting in the DPO model. This process involves LoRA fine-tuning with a learning rate of $7e-6$ trained for 3 epochs, and the maximum input length during training is set to 10000 tokens. As in the previous stage, we leverage DeepSpeed ZeRO-3 and FlashAttention2 for efficient training, with a total batch size of 32 and BF16 precision.

In the Query Aware Policy Optimization stage, we focus on a curated set of challenging questions that remained unresolved after four sampling trials. Through RL, the DPO model is further optimized to produce the final SmartSearch model. The training is

conducted with a learning rate of $1e-6$, where each sample undergoes 8 rollouts to explore different trajectories. The total training batch size is 64, with a PPO mini-batch size of 16. We set the maximum output length to 8192 tokens and limit the number of tool calls to 5 during each rollout.

During the Inference stage, we set the maximum output length to 16384 tokens and limit the number of tool calls to 10.

C Case Study

To better demonstrate the effectiveness of SmartSearch, as well as the process reward and query refinement mechanisms, we conducted a case study.

As shown in Table 4, when handling a user’s question, SmartSearch first comprehends the question and performs preliminary planning to form an accurate search intent. Based on this intent, the model formulate a precise search query and successfully retrieves the desired information. It is also noteworthy that our model effectively utilizes its internal knowledge, generating only necessary

search queries, which further enhances the search efficiency of SmartSearch.

Table 5 illustrates how process reward mechanism provides both numerical scores and textual explanations for each search query within the model’s output. The scores accurately reflect the quality of the search queries, while the explanations offer detailed feedback. For low-quality queries, the explanations clearly identify the reasons for their subpar performance, providing crucial guidance for subsequent refinement.

As presented in Table 6, query refinement mechanism refines the low-quality queries based on the explanations provided by the process reward and regenerates the subsequent steps from the refined query. The initial query fails to retrieve the expected information, causing the entire trajectory to deviate from the correct path. In contrast, the refined query successfully retrieves the desired information and leads to the correct answer. The comparison between the two trajectories helps the model better concentrate its optimization on search queries, highlighting the effectiveness of this mechanism.

Table 4: Example of SmartSearch output, including the question, golden answer, and model output.

Example of SmartSearch Output
<p>Question: Douglas D. Scott is an American archaeologist most notable for his work at the site of a battle that occurred in what year?</p> <p>Golden Answer: 1876</p> <p>Model Output: <think> Okay, so I need to figure out the year of the battle that Douglas D. Scott worked on. Let me start by recalling what I know about Douglas D. Scott. He's an American archaeologist, so maybe he's known for excavating a significant site. The question mentions a battle site, so I should look for his notable archaeological sites related to battles. First, I'll try to search for Douglas D. Scott's main archaeological sites. Maybe he's famous for a specific battle site. Let me start with a general search query. </think> <search> Douglas D. Scott notable archaeological sites </search> <result> result: "Douglas D. Scott" by Ronald K. Wetherington and Frances Levine, pp 134–152, University of Oklahoma Press, Norman. Douglas D. Scott Douglas D. Scott is an American archaeologist most notable for his work at the Little Bighorn in the mid-1980s. Working with Richard Fox, Melissa Connor, Doug Harmon, and staff and volunteers from the National Park Service, Scott worked to sketch out a field methodology that has enabled archaeologists to systematically investigate battlefields. This work is internationally recognized as constituting a great step forward in our ability to interpret battlefields archaeologically, regardless of the extent of the historical record. At the Little Bighorn, the ... </result> <think> Hmm, looking at the search results, the first entry mentions Douglas D. Scott's work at the Little Bighorn in the mid-1980s. The other entries don't seem to mention a battle site. The Little Bighorn is a well-known battle, so that's probably the one. The question asks for the year of the battle. The Little Bighorn was a major battle during the American Indian Wars, specifically between the Lakota, Northern Cheyenne, and Arapaho tribes against the U.S. Army. The Battle of the Little Bighorn took place on June 25, 1876. Let me confirm the exact year. The search result says 'mid-1980s' for his work at Little Bighorn. Since the battle itself was in 1876, the year would be 1876. I don't see any conflicting information here. The other sites mentioned are not battles but other archaeological sites. So the answer should be 1876. </think> <answer> The final answer is 1876. </answer></p>

Table 5: Example of process reward, including the question, golden answer, model output, and process reward for each query within the model output.

Example of Process Reward
<p>Question: An Annapolis Story stars which American stage, film, and television actor born on February 15, 1914?</p> <p>Golden Answer: Kevin McCarthy</p> <p>Model Output: <think> To solve this question, I need to find out which American stage, film, and television actor born on February 15, 1914, stars in the movie “An Annapolis Story”. Step 1: Identify the actor who stars in “An Annapolis Story”. Step 2: Determine which of that actor’s birthdate matches February 15, 1914. Let’s start with Step 1: Step 1: Identify the actor who stars in “An Annapolis Story”. I’ll use a query to search for information about the stars of “An Annapolis Story”: </think> <search> stars in “An Annapolis Story” </search> <result> result:“An Annapolis Story” An Annapolis Story An Annapolis Story (alternative titles: The Blue and Gold and Navy Air Patrol) is a 1955 American drama film directed by Don Siegel and starring John Derek, Diana Lynn and Kevin McCarthy.A product of the newly formed Allied Artists company, “An Annapolis Story”, despite having the “Siegel Touch”, suffered from its low budget. Brothers Tony (John Derek) and Jim Scott (Kevin McCarthy) enroll as midshipmen at the United States Naval Academy in Annapolis. Jim, the older one, looks after the more impulsive Tony and helps him pass a difficult test so he can play football in the ... </result> <think> From the search results, I can see that “An Annapolis Story” is a 1955 American drama film starring John Derek, Diana Lynn, and Kevin McCarthy. Step 2: Determine which of that actor’s birthdate matches February 15, 1914. I’ll use a query to search for information about the birthdate of Kevin McCarthy: </think> <search> birthdate of Kevin McCarthy </search> <result> result: “Kevin McCarthy (California politician)” Kevin McCarthy (California politician) Kevin Owen McCarthy (born January 26, 1965) is an American politician serving as the House Majority Leader since 2014 and U.S. Representative for California’s 23rd congressional district since 2007. The 23rd district, numbered as the 22nd district from 2007 to 2013, is based in Bakersfield and includes large sections of Kern County and Tulare County as well as part of the Quartz Hill neighborhood in northwest Los Angeles County. A member of the Republican Party, he was formerly chairman of the California Young Republicans and the Young Republican National Federation. McCarthy worked as district director for ... </result> <think> From the search results, I can see that Kevin McCarthy was born on January 26, 1965. Therefore, the American stage, film, and television actor born on February 15, 1914, who stars in “An Annapolis Story” is not Kevin McCarthy. However, the question might be referring to another actor. ... </think></p> <p>Process Reward: Score1: 1 Explanation1: The search intent is necessary, and the search results included the cast list for the movie An Annapolis Story. Score2: 0 Explanation2: The search intent is necessary, but the search results did not include the birth date of actor Kevin McCarthy. Instead, they contained information about politician Kevin McCarthy. Therefore, the score is 0.</p>

Table 6: Example of query refinement, including the question, golden answer, model output, refined query, and the regenerated subsequent steps.

Example of Query Refinement
<p>Question: An Annapolis Story stars which American stage, film, and television actor born on February 15, 1914?</p>
<p>Golden Answer: Kevin McCarthy</p>
<p>Original Model Output: <think> To solve this question, I need to find out which American stage, film, and television actor born on February 15, 1914, stars in the movie “An Annapolis Story”. Step 1: Identify the actor who stars in “An Annapolis Story”. Step 2: Determine which of that actor’s birthdate matches February 15, 1914. Let’s start with Step 1: Step 1: Identify the actor who stars in “An Annapolis Story”. I’ll use a query to search for information about the stars of “An Annapolis Story”: </think> <search> stars in “An Annapolis Story” </search> <result> result: “An Annapolis Story” An Annapolis Story An Annapolis Story (alternative titles: The Blue and Gold and Navy Air Patrol) is a 1955 American drama film directed by Don Siegel and starring John Derek, Diana Lynn and Kevin McCarthy.A product of the newly formed Allied Artists company, “An Annapolis Story”, despite having the “Siegel Touch”, suffered from its low budget. Brothers Tony (John Derek) and Jim Scott (Kevin McCarthy) enroll as midshipmen at the United States Naval Academy in Annapolis. Jim, the older one, looks after the more impulsive Tony and helps him pass a difficult test so he can play football in the ... </result> <think> From the search results, I can see that “An Annapolis Story” is a 1955 American drama film starring John Derek, Diana Lynn, and Kevin McCarthy. Step 2: Determine which of that actor’s birthdate matches February 15, 1914. I’ll use a query to search for information about the birthdate of Kevin McCarthy: </think> <search> birthdate of Kevin McCarthy </search> <result> result: “Kevin McCarthy (California politician)” Kevin McCarthy (California politician) Kevin Owen McCarthy (born January 26, 1965) is an American politician serving as the House Majority Leader since 2014 and U.S. Representative for California’s 23rd congressional district since 2007. The 23rd district, numbered as the 22nd district from 2007 to 2013, is based in Bakersfield and includes large sections of Kern County and Tulare County as well as part of the Quartz Hill neighborhood in northwest Los Angeles County. A member of the Republican Party, he was formerly chairman of the California Young Republicans and the Young Republican National Federation. McCarthy worked as district director for ... </result> <think> From the search results, I can see that Kevin McCarthy was born on January 26, 1965. Therefore, the American stage, film, and television actor born on February 15, 1914, who stars in “An Annapolis Story” is not Kevin McCarthy. However, the question might be referring to another actor. ... </think> Query Refinement <search> birthdate of Kevin McCarthy </search> → <search> birthdate of Actor Kevin McCarthy </search> Regenerate Subsequent Steps <search> birthdate of Actor Kevin McCarthy </search> <result> result: “Kevin McCarthy (actor)” Kevin McCarthy (actor) Kevin McCarthy (February 15, 1914 – September 11, 2010) was an American actor who gave over 200 television and film performances. He is best remembered for portraying the male lead in the horror science fiction film “Invasion of the Body Snatchers” (1956). Following several television guest roles, McCarthy gave his first credited film performance in “Death of a Salesman” (1951), portraying Biff Loman to Fredric March’s Willy Loman. The role earned him a Golden Globe Award and a nomination for the Academy Award for Best Supporting Actor. McCarthy was born in Seattle, Washington, the son of Roy ... </result> <think> The search results confirm that Kevin McCarthy was an American actor born on February 15, 1914. He starred in “An Annapolis Story” (1955) and had a long career in film, television, and stage. The question asks for an actor who starred in “An Annapolis Story” and was born on February 15, 1914. Kevin McCarthy fits all these criteria. </think> <answer> Kevin McCarthy </answer></p>