

ZEROSEARCH: Incentivize the Search Capability of LLMs without Searching

Hao Sun, Zile Qiao*, Jiayan Guo*, Xuanbo Fan, Yingyan Hou
Yong Jiang, Pengjun Xie, Fei Huang, Yan Zhang*

Tongyi Lab , Alibaba Group
 [Homepage](#)  [Model](#)  [Datasets](#)  [Code](#)

Abstract

Effective information searching is essential for enhancing the reasoning and generation capabilities of large language models (LLMs). Recent research has explored using reinforcement learning (RL) to improve LLMs’ search capabilities by interacting with live search engines in real-world environments. While these approaches show promising results, they face two major challenges: (1) **Uncontrolled Document Quality**: The quality of documents returned by search engines is often unpredictable, introducing noise and instability into the training process. (2) **Prohibitively High API Costs**: RL training requires frequent rollouts, potentially involving hundreds of thousands of search requests, which incur substantial API expenses and severely constrain scalability. To address these challenges, we introduce **ZEROSEARCH**, a reinforcement learning framework that incentivizes the search capabilities of LLMs without interacting with real search engines. Our approach begins with lightweight supervised fine-tuning to transform the LLM into a retrieval module capable of generating both relevant and noisy documents in response to a query. During RL training, we employ a curriculum-based rollout strategy that incrementally degrades the quality of generated documents, progressively eliciting the model’s reasoning ability by exposing it to increasingly challenging retrieval scenarios. Extensive experiments demonstrate that ZEROSEARCH effectively incentivizes the search capabilities of LLMs using a 3B LLM as the retrieval module. Remarkably, a 7B retrieval module achieves comparable performance to the real search engine, while a 14B retrieval module even surpasses it. Furthermore, it generalizes well across both base and instruction-tuned models of various parameter sizes and is compatible with a wide range of RL algorithms.

1 Introduction

Large Language Models (LLMs) [36, 3, 45] have demonstrated remarkable performance across a wide range of downstream tasks, including mathematical reasoning, question answering, and code generation [38, 39, 11, 22]. However, the knowledge encoded in these models is inherently static, constrained by the scope of the data encountered during pretraining. As a result, LLMs remain prone to generating hallucinated content or outdated information [13, 34, 44], which undermines their reliability in practical applications. Therefore, it is essential to enable LLMs to access external sources of information to produce more accurate and grounded responses.

One widely adopted approach to addressing this issue is Retrieval-Augmented Generation (RAG), which incorporates external knowledge into the generation pipeline [29, 33, 30, 6, 2, 27]. Early

*Corresponding Author.

work in this area focused on prompt-based strategies that guide LLMs through query generation, query decomposition, and multi-turn information retrieval [43, 28, 42, 16, 33, 23]. While effective, these strategies often require meticulous prompt engineering and place high demands on the model’s reasoning capabilities. To improve the efficiency, subsequent research explored supervised fine-tuning (SFT) to enhance the performance of smaller LLMs [1, 24, 12]. Further advances have focused on test-time scaling techniques [25, 15, 46, 14], such as Monte Carlo Tree Search (MCTS), which dynamically expands the search space during inference. Although promising, such methods incur significant computational overhead, posing challenges for practical deployment.

Recently, reinforcement learning (RL) has emerged as a promising strategy for further improving LLM performance by enhancing their reasoning and decision-making capabilities [7, 9]. Notably, RL-based models such as OpenAI-o1 and DeepSeek-R1 have demonstrated substantial gains in logical inference and iterative reasoning—achieved purely through reward-driven learning, without relying on explicit step-by-step supervision [20]. Within this paradigm, several studies have explored using RL to train policy models that can more effectively search for relevant information. Representative examples include Search-R1 [17], R1-Searcher [35], and ReSearch [19]. Notably, DeepResearcher [47] introduces live interaction with commercial search engines such as Google, allowing models to train in an environment that closely mirrors real-world web search. Despite these advancements, integrating RL with real-world search scenarios presents significant challenges: (1) **Uncontrolled Document Quality**: The quality of documents retrieved from live search engines is often unpredictable, introducing noise and instability into the training process. (2) **Prohibitively High API Costs**: RL training requires frequent rollouts, potentially involving hundreds of thousands of API calls, which incur substantial financial costs and severely limit scalability.

To address these challenges, we propose **ZEROSEARCH**—a reinforcement learning framework that enables LLMs to learn search strategies without interacting with real search engines. Our key insight is that **LLMs have acquired extensive world knowledge during large-scale pretraining and are capable of generating relevant documents given a search query** [43]. The primary difference between a real search engine and a simulation LLM lies in the textual style of the returned content. However, with lightweight supervised fine-tuning, even relatively small LLMs can effectively simulate the behavior of real search engines. In addition to eliminating API costs, an important advantage of using LLMs for document generation is the ability to control document quality. Specifically, during supervised fine-tuning, documents that lead to correct or incorrect answers are distinguished through prompt design, enabling the simulation LLM to learn to generate either relevant or noisy documents simply by adjusting a few words in the prompt. Building on this, **we introduce a curriculum rollout mechanism during training**, in which the quality of the generated documents is gradually degraded over time to simulate increasingly challenging retrieval scenarios. This allows the policy model to first learn basic output formats and task requirements before progressively adapting to more challenging and noisy retrieval scenarios. More importantly, ZEROSEARCH exhibits strong scalability: increasing the number of GPUs significantly accelerates the generation throughput of the simulation LLM, thereby enabling efficient large-scale rollout. Empirical results show that **even a 3B LLM used as the simulated search engine can effectively incentivize the policy model’s search capabilities**. A 7B retrieval module achieves performance comparable to Google Search, while a 14B retrieval module even surpasses it. ZEROSEARCH is compatible with both base and instruction-tuned models of various parameter sizes, removing the need for separate supervised warm-up stages. Moreover, it integrates seamlessly with widely used RL algorithms, including Proximal Policy Optimization (PPO) [31], Group Relative Policy Optimization (GRPO) [32, 7], and Reinforce++ [10].

Our contributions can be summarized as follows:

- We propose ZEROSEARCH, a novel reinforcement learning framework that incentivizes the search capability of LLMs without interacting with real search engines.
- Through supervised fine-tuning, we transform the LLM into a retrieval module capable of generating both relevant and noisy documents in response to a query. We further introduce a curriculum rollout mechanism to progressively elicit the model’s reasoning ability by exposing it to increasingly challenging retrieval scenarios.
- We conduct extensive experiments on both in-domain and out-of-domain datasets. Results show that ZEROSEARCH outperforms real search engine-based models while incurring zero API cost. Moreover, it generalizes well across both base and instruction-tuned LLMs of various parameter sizes and supports different reinforcement learning algorithms.

2 Related Work

2.1 Retrieval-Augmented Generation

Retrieval-augmented generation (RAG) enhances generation performance by integrating relevant external knowledge into the generation pipeline. Early research primarily adopted prompt-based approaches, guiding LLMs through processes such as query generation, query decomposition, and multi-turn information retrieval [43, 28, 42, 16, 33, 23]. Despite their effectiveness, these methods often require intricate prompt engineering and impose substantial demands on the model’s reasoning capabilities. To improve efficiency and reduce dependency on strong black-box LLMs, subsequent work has proposed supervised fine-tuning strategies for smaller LLMs. For instance, Self-RAG [1] employs a self-reflection mechanism, iteratively refining model outputs through predicted reflection tokens. RetroLLM [24] integrates retrieval and generation by enabling the model to directly generate fine-grained evidence from the corpus via constrained decoding. Recent advances also include test-time scaling techniques [25, 15, 46, 14], notably Monte Carlo Tree Search (MCTS), which dynamically expands the search space during inference. For example, RAG-star [14] integrates retrieved information into a tree-based reasoning process, while AirRAG [5] employs MCTS to activate intrinsic reasoning capabilities and expand the solution space. Despite promising results, these approaches introduce significant computational overhead, limiting their practical applicability.

2.2 Learning to Search through Reinforcement Learning

Recently, reinforcement learning (RL) has emerged as a promising paradigm for enhancing the reasoning capabilities of LLMs [7, 9]. Notable RL-based models such as OpenAI-o1 and DeepSeek-R1 have demonstrated remarkable capabilities in logical inference and iterative reasoning, purely driven by reward signals without explicit step-by-step supervision [20]. Several studies have also explored RL techniques specifically designed to train models for effective information retrieval. For instance, Search-R1 [17] employs reinforcement learning to autonomously generate multiple search queries during step-by-step reasoning. Similarly, R1-Searcher [35] proposes a two-stage, outcome-based RL method aimed at enhancing search capabilities. ReSearch [19] leverages RL to teach models to reason through searches, entirely without supervision on intermediate reasoning steps. However, these methods usually employ static, local textual corpora such as Wikipedia and fail to capture the complexities of real-world interaction. To bridge this gap, DeepResearcher [47] introduces direct interaction with commercial search engines such as Google, allowing training environments that closely approximate real-world search scenarios. While achieving superior performance, these real-time retrieval methods face significant challenges, including unpredictable document quality, prohibitive high API costs that adversely affect system scalability. To address these limitations, we propose ZEROSEARCH, a method leveraging an LLM to simulate real-time search, effectively removing dependence on costly, rate-limited real search APIs. Through lightweight supervised fine-tuning, ZEROSEARCH allows explicit control over document quality and implements a curriculum rollout mechanism that enhances training stability and robustness.

3 ZEROSEARCH

In this section, we first formalize the reinforcement learning objective without a search engine. We then detail the design of ZEROSEARCH, covering the training template, search simulation tuning, curriculum-based rollout strategy, reward design, and training algorithms.

3.1 Reinforcement Learning without a Search Engine

We propose a reinforcement learning framework that eliminates the need for a real search engine by leveraging an LLM to simulate the search engine. The optimization objective is formulated as:

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot | x; \pi_{\psi})} [r_{\phi}(x, y)] - \beta \text{D}_{\text{KL}}[\pi_{\theta}(y | x; \pi_{\psi}) \parallel \pi_{\text{ref}}(y | x; \pi_{\psi})],$$

where π_{θ} is the policy model to be optimized, π_{ref} is the reference model, and r_{ϕ} denotes the reward function. π_{ψ} represents the simulation LLM, whose parameters remain fixed throughout training.

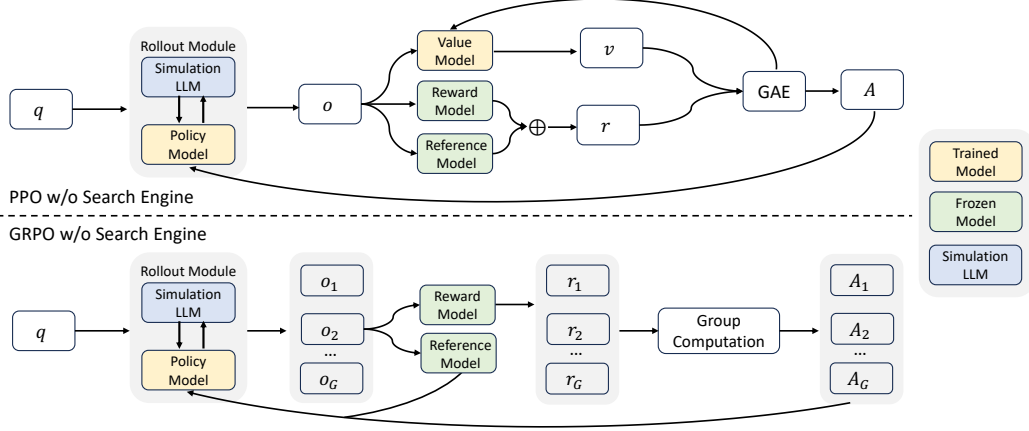


Figure 1: Demonstration of PPO and GRPO training without the search engine.

3.2 Training Template

Answer the given question. You must conduct reasoning inside `<think>` and `</think>` first every time you get new information. After reasoning, if you find you lack some knowledge, you can call a search engine by `<search>` query `</search>`, and it will return the top searched results between `<information>` and `</information>`. You can search as many times as you want. If you find no further external knowledge needed, you can directly provide the answer inside `<answer>` and `</answer>` without detailed illustrations. For example, `<answer>` Beijing `</answer>`. Question:

Table 1: Training template. The question is appended at the end during training and inference.

In ZEROSEARCH, rather than relying on supervised fine-tuning for generation, we follow [17] and apply a multi-turn interaction template that guides the policy model through iterative reasoning and information retrieval until a final answer is reached.

As illustrated in Table 1, the interaction is divided into three distinct stages: First, the model articulates its internal reasoning within the `<think>...</think>` tag. Second, if additional evidence is needed, it issues a search query within the `<search>...</search>` tag. Finally, once sufficient information has been retrieved, the model provides its answer in the `<answer>...</answer>` tag. This explicit separation of reasoning, searching, and answering enforces a structured decision-making process, enhancing the model’s transparency and reliability.

3.3 Search Simulation Tuning

During rollout, we use the LLM to simulate a real search engine by generating documents in response to queries. A straightforward approach is to directly prompt the LLM to generate documents. However, this often results in a noticeable style gap compared to outputs from real search engines.

To bridge this gap, we propose a lightweight supervised fine-tuning (SFT) procedure. Specifically, we first collect interaction trajectories by prompting the LLM to engage with a real search engine in a multi-turn manner until a final answer is reached. Trajectories that yield correct answers are labeled as positive, indicating that the retrieved documents were useful, while those leading to incorrect answers are labeled as negative, indicating noisy retrievals.

Then, we extract query-document pairs from both positive and negative trajectories and perform lightweight SFT to improve the LLM’s ability to simulate real search engines. As shown in Table 2, the distinction between useful and noisy retrievals is achieved by adjusting a few words in the prompt. Besides, we also incorporate the input question and its corresponding answer into the prompt to broaden the knowledge boundary of LLMs. After fine-tuning, the LLM is capable of generating both useful and noisy documents, enabling dynamic document quality control during rollout.

3.4 Rollout with Curriculum Search Simulation

You are the Google search engine.
 Given a query, you need to generate five [useful / noisy] documents for the query.
 The user is trying to answer the question: [question] whose answer is [ground truth].
 Each document should contain about 30 words, and these documents should contain [useful / noisy] information.
 Query: [query]
 [Useful / Noisy] Output:

Table 2: Template for Search Simulation. The **useful** and **noisy** keywords are used to control the quality of the generated documents. The input question and its ground-truth answer are also included in the prompt to help extend the simulation LLM’s knowledge coverage.

During rollout, the policy model performs interactive reasoning and generates search queries, which are fed into the simulation LLM to produce corresponding documents. To gradually increase the difficulty of training, we introduce a curriculum learning-based rollout mechanism, where the quality of the generated documents is progressively degraded over time. This is controlled by a probability function p_i that governs the likelihood of generating noisy documents at step i :

$$p_i = p_s + \frac{b^{i/m} - 1}{b - 1} (p_e - p_s) \quad (1)$$

Here, p_s and p_e represent the initial and final noise probabilities, i and m denote the current and total number of training steps, and b is the exponential base, with a default value of 4. As training progresses, the ratio i/m increases, leading to a higher p_i value—*i.e.*, a greater chance of producing noisy documents. This allows the policy model to first learn basic output structures and task requirements, before gradually adapting to more challenging and noisy retrieval scenarios.

3.5 Reward Design

The reward signal serves as the primary supervision in the reinforcement learning process. In this work, we adopt a rule-based reward function that focuses solely on answer accuracy. During preliminary experiments, we observed that using exact match (EM) as the reward metric often led to reward hacking: the policy model tended to produce excessively long answers to increase the chance of including the correct answer. To mitigate this issue, we adopt an F1 score-based reward, which balances precision and recall, and is calculated as:

$$r_\phi(x, y) = \frac{2 \times IN}{PN + RN},$$

where IN denotes the number of overlapping words between the prediction and the ground truth, PN is the number of words in the prediction, and RN is the number of words in the ground truth. We do not incorporate an additional reward for output format, as we observe that the model consistently produces well-formed responses without explicit supervision.

3.6 Training Algorithm

Our approach is compatible with a wide range of reinforcement learning algorithms, including Proximal Policy Optimization (PPO) [31], Group Relative Policy Optimization (GRPO) [32, 7], and Reinforce++ [10], each offering distinct advantages for optimizing retrieval-augmented reasoning.

In ZEROSEARCH, the rollout sequence comprises both tokens generated by the policy model and document tokens returned by the simulation LLM. Applying the same optimization procedure uniformly across both types of tokens can lead to training instability, as the retrieved content is externally generated and not directly controlled by the policy model.

To mitigate this, we introduce a loss masking mechanism for retrieved tokens, ensuring that gradients are only computed with respect to the model’s own outputs. This strategy stabilizes the RL training process while preserving the effectiveness of retrieval-augmented generation.

4 Main Results

4.1 Datasets and Evaluation Metrics

We evaluate ZEROSEARCH on a diverse set of question answering benchmarks: (1) **Single-Hop Question Answering**, including NQ [21], TriviaQA [18], and PopQA [26]. (2) **Multi-Hop Question Answering**, including HotpotQA [41], 2WikiMultiHopQA [8], Musique [37], and Bamboogle [28].

We follow [17] and adopt Exact Match (EM) as our evaluation metric. A prediction is deemed correct if its normalized form exactly matches any of the normalized ground-truth answers.

4.2 Baselines

To evaluate the effectiveness of ZEROSEARCH, we compare our method with the following baselines. (1) **Vanilla Prompting Methods**: This category includes direct prompting, Chain-of-Thought (CoT), and standard Retrieval-Augmented Generation (RAG). (2) **Advanced RAG Methods**: We consider RAgent [23] and Search-o1 [23], which iteratively search for relevant information. (3) **RL Tuning Methods**: This category includes R1 and Search-R1 [17]. In R1, the policy model is trained to perform in-depth reasoning based solely on its internal knowledge. In contrast, Search-R1 enables the policy model to interact with a real search engine multiple times during inference.

To ensure a fair comparison, we adopt the F1 score as the reward metric across all RL methods. Notably, among RL-based search baselines, **we compare only with Search-R1, as it avoids complex reward design, data selection, or elaborate training pipelines**. This setting allows for a direct and equitable comparison between the real search engine and our simulated search engine.

4.3 Experimental Setup

We conduct experiments using three model families: Qwen-2.5-7B (Base/Instruct) and Qwen-2.5-3B (Base/Instruct) [40], as well as LLaMA-3.2-3B (Base/Instruct) [4]. To simulate real-world retrieval scenarios, we utilize Google Web Search via the SerpAPI² as the external search engine. The number of retrieved documents is fixed at five across all methods to ensure a fair comparison.

For datasets, following the setup in [17], we merge the training sets of NQ and HotpotQA to create a unified dataset for all fine-tuning-based approaches. Evaluation is conducted on seven datasets to assess both in-domain and out-of-domain performance. For prompt-based baselines, we use Instruct models, as Base models typically struggle to follow task instructions. For RL-based methods, we evaluate both Base and Instruct variants to assess generality across model types.

To train the simulation LLM, we conduct a lightweight SFT using Qwen-2.5-3B, Qwen-2.5-7B, and Qwen-2.5-14B as the backbones. The learning rate is set to be $1e-6$. To train ZEROSEARCH, we adopt two reinforcement learning algorithms: GRPO and PPO. In the GRPO setting, the policy LLM is trained with a learning rate of $1e-6$, and five responses are sampled per prompt. In the PPO setting, the policy LLM is trained with a learning rate of $1e-6$, while the value model is trained with a separate learning rate of $1e-5$. We apply Generalized Advantage Estimation (GAE) with hyperparameters $\lambda = 1$ and $\gamma = 1$. Unless otherwise specified, GRPO is used as the default reinforcement learning algorithm, and Qwen-2.5-14B is used as the default simulation LLM throughout all experiments.

4.4 Performance

Table 3 presents a comparison between ZEROSEARCH and several baseline methods across seven datasets. Based on the results, several key observations can be drawn:

ZEROSEARCH consistently outperforms all baseline methods. This performance advantage holds for both in-domain datasets (*i.e.*, NQ and HotpotQA) and out-of-domain datasets (*i.e.*, TriviaQA, PopQA, 2WikiMultiHopQA, Musique, and Bamboogle), demonstrating the robustness of our method.

ZEROSEARCH surpasses methods that rely on real search engines. Compared to Search-R1, which utilizes the real search engine, ZEROSEARCH achieves better performance, highlighting its potential as an effective alternative to real search engines in large-scale reinforcement learning.

²<https://serpapi.com/>

Method	Single-Hop QA			Multi-Hop QA				
	NQ	TriviaQA	PopQA	HotpotQA	2Wiki	Musique	Bamboogle	Avg.
<i>Qwen-2.5-7B-Base/Instruct</i>								
Direct Answer	11.60	35.60	1.20	16.40	22.20	4.80	14.40	15.17
CoT	12.80	35.60	3.80	16.20	22.60	6.60	24.00	17.37
RAG	27.40	58.20	17.80	25.80	23.20	9.40	16.80	25.51
RA-Agent	21.20	40.20	8.80	19.60	19.60	7.60	28.00	20.71
Search-o1	19.40	40.60	11.40	17.00	27.00	8.60	30.40	22.06
R1-base	25.15	43.18	22.29	21.02	28.46	9.76	24.80	24.95
R1-instruct	25.25	42.68	27.81	20.45	26.83	8.33	27.05	25.49
Search-R1-base	41.51	60.53	51.02	32.25	36.31	16.39	28.00	38.00
Search-R1-inst	41.46	62.17	49.80	34.55	34.22	19.43	33.06	39.24
ZEROSearch-base	41.84	63.54	51.72	30.30	40.33	12.25	30.25	38.61
ZEROSearch-inst	43.24	61.81	51.52	29.21	43.12	19.72	35.20	40.54
<i>Qwen-2.5-3B-Base/Instruct</i>								
Direct Answer	12.40	30.60	5.60	16.00	19.20	4.40	16.80	15.00
CoT	15.00	33.60	3.60	16.20	18.00	3.60	12.80	14.69
RAG	31.60	58.00	15.20	24.20	23.20	8.20	15.20	25.09
RA-Agent	15.20	28.40	6.60	12.60	16.60	2.60	13.60	13.66
Search-o1	16.60	31.00	8.20	14.80	22.40	5.20	22.40	17.23
R1-base	21.46	34.02	16.29	22.67	28.57	5.51	7.20	19.39
R1-instruct	17.41	33.13	18.40	18.90	25.66	4.69	15.32	19.07
Search-R1-base	41.13	61.22	40.73	27.64	31.97	12.17	12.40	32.47
Search-R1-inst	40.82	59.11	42.76	30.83	31.10	8.37	13.01	32.28
ZEROSearch-base	41.21	61.49	43.99	31.02	33.20	12.58	14.29	33.97
ZEROSearch-inst	66.19	57.29	23.39	24.44	30.75	9.39	19.67	33.02
<i>LLaMA-3.2-3B-Base/Instruct</i>								
Direct Answer	16.20	29.60	7.40	12.60	9.20	2.00	8.00	12.14
CoT	26.20	44.40	2.80	16.00	10.20	5.80	21.60	18.14
RAG	30.00	57.60	26.40	23.40	17.60	9.60	11.20	25.11
RA-Agent	22.40	36.20	11.40	16.60	21.00	5.60	26.40	19.94
Search-o1	24.20	48.40	8.80	19.40	17.40	6.00	32.00	22.31
R1-base	26.63	42.31	20.16	22.56	28.31	7.36	15.45	23.25
R1-instruct	32.32	52.43	25.98	19.51	29.21	10.34	21.95	27.39
Search-R1-base	40.85	61.68	44.47	28.25	32.66	13.74	12.10	33.39
Search-R1-inst	43.72	56.44	45.49	22.45	23.37	9.27	18.70	31.35
ZEROSearch-base	40.28	61.87	49.80	30.57	35.05	14.75	14.52	35.26
ZEROSearch-inst	40.52	60.41	48.99	20.98	27.40	6.35	13.71	31.19

Table 3: Main results using different LLMs as the backbone. The best performance is set in bold.

ZEROSearch demonstrates strong generalizability. Across different model families, parameter sizes, and types (*i.e.*, base or instruction-tuned), ZEROSearch consistently outperforms baselines. Moreover, its performance further improves with larger models, highlighting its scalability.

5 Further Analysis

5.1 Compare ZEROSearch with Real Search Engine

We compare the reward curves of ZEROSearch and Search-R1 (using a real search engine) on LLaMA-3.2-3B, as shown in Figures 2a and 2b. Several key observations can be made:

The overall reward trends are similar across both methods. As training progresses, the reward scores of both ZEROSearch and Search-R1 steadily increase, indicating that the policy models in both settings effectively learn to interact with search engines and produce correct answers.

ZEROSearch achieves a more stable and smoother learning curve. As shown in Figure 2b, ZEROSearch initially lags behind Search-R1 but eventually surpasses it with much less fluctuation, thanks to the curriculum mechanism that helps the model gradually master search tool usage.

ZEROSearch generalizes well across both base and instruction-tuned models. Under both model types, ZEROSearch steadily improves reward performance, underscoring its generalizability.

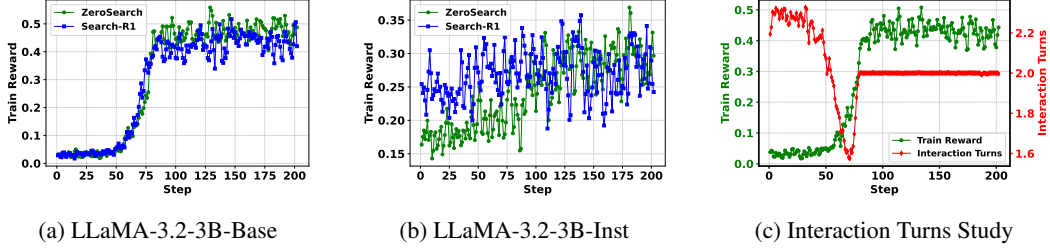


Figure 2: (a-b): Reward curve comparison between ZEROSEARCH and Search-R1 using LLaMA-3.2-3B. (c): Interaction turns and reward progression during training of LLaMA-3.2-3B-base.

Search Engine	NQ	TriviaQA	PopQA	HotpotQA	2Wiki	Musique	Bamboogle	Avg.
Prompt-3B	35.77	52.74	41.62	25.76	26.02	8.47	10.57	28.71
Prompt-7B	39.71	57.93	38.24	29.47	27.64	7.93	7.26	29.74
Prompt-14B	40.40	58.82	37.78	26.26	29.01	10.14	15.45	31.12
SFT-3B	42.03	59.68	44.22	29.18	30.24	10.41	11.29	32.44
SFT-7B	41.70	61.18	46.46	30.66	28.98	11.76	10.66	33.06
SFT-14B	41.21	61.49	43.99	31.02	33.20	12.58	14.29	33.97
Google	41.13	61.22	40.73	27.64	31.97	12.17	12.40	32.47

Table 4: Performance of simulated search engines using different LLM configurations. We compare prompt-based and fine-tuned simulation LLMs (3B to 14B) against Google Search.

5.2 Choice of Simulation LLMs

In this section, we investigate how different simulation engine configurations affect performance, including prompt-based and fine-tuned LLMs ranging from 3B to 14B parameters. Based on the results in Table 4, we make the following observations:

First, the fine-tuned 7B simulation engine (SFT-7B) achieves performance comparable to that of Google Search, while the 14B variant (SFT-14B) even surpasses it. This demonstrates the feasibility of using a well-trained LLM as a substitute for real search engines in reinforcement learning setups.

Second, fine-tuned simulation engines significantly outperform prompt-based ones. Although prompt-based methods are explicitly guided to mimic the response style of a real search engine, a substantial distribution gap remains, leading to inferior performance.

Third, performance improves consistently with increasing model size. Larger simulation LLMs not only exhibit stronger simulation capabilities but also more accurately distinguish between relevant and irrelevant documents, thereby enabling more effective curriculum learning during training.

5.3 Interaction Turns Study

In this section, we analyze the training dynamics of ZEROSEARCH by examining both reward progression and the number of interaction turns throughout the training, using the LLaMA3.2-3B-Base model. The results are shown in Figure 2c.

During the early phase of training, the number of interaction turns drops sharply, while the reward increases slowly. This is primarily because the policy model initially lacks knowledge of how to properly invoke the search engine, resulting in redundant interactions. However, it quickly learns the correct format and begins to eliminate unnecessary steps effectively.

As training progresses, both the number of interaction turns and the reward curve increase sharply and then stabilize. This is primarily because the policy model becomes capable of effectively retrieving relevant documents and ultimately achieving correct answers, resulting in higher rewards. Notably, although the reward appears stable in the later stages of training, the underlying task difficulty continues to rise due to the curriculum mechanism. Consequently, the policy must continuously refine its strategy and improve its reasoning ability to maintain consistent performance.

Method	NQ	TriviaQA	PopQA	HotpotQA	2Wiki	Musique	Bamboogle	Avg.
<i>Qwen-2.5-3B-Base</i>								
PPO	39.92	60.73	46.25	29.63	34.42	11.90	16.53	34.20
GRPO	41.21	61.49	43.99	31.02	33.20	12.58	14.29	33.97
<i>LLaMA-3.2-3B-Base</i>								
PPO	39.92	62.15	43.61	29.12	35.44	12.98	15.32	34.08
GRPO	40.28	61.87	49.80	30.57	35.05	14.75	14.52	35.26

Table 5: Performance of ZEROSEARCH under different RL algorithms. We compare PPO and GRPO using the Qwen2.5-3B-Base and LLaMA-3.2-3B-Base models.

Method	NQ	TriviaQA	PopQA	HotpotQA	2Wiki	Musique	Bamboogle	Avg.
<i>Qwen-2.5-3B-Base</i>								
Easy → Hard	41.21	61.49	43.99	31.02	33.20	12.58	14.29	33.97
Hard → Easy	39.59	58.78	44.31	32.32	30.67	12.55	12.00	32.89
<i>Qwen-2.5-3B-Instruct</i>								
Easy → Hard	66.19	57.29	23.39	24.44	30.75	9.39	19.67	33.02
Hard → Easy	40.16	57.83	24.69	28.72	27.55	8.96	8.94	28.12

Table 6: Reverse Curriculum Study. We compare the performance of standard and reverse curriculum rollout settings using the Qwen-2.5-3B-Base and Qwen-2.5-3B-Instruct models.

5.4 Different RL algorithms: PPO vs. GRPO

In this section, we evaluate the performance of two widely adopted RL training algorithms, PPO and GRPO, within the ZEROSEARCH framework, using the Qwen2.5-3B-Base and LLaMA-3.2-3B-Base models. The results of this comparison are presented in Table 5.

As observed, both GRPO and PPO successfully incentivize the search capabilities within our framework, demonstrating the versatility of our approach. Among them, GRPO exhibits more stable performance across both models, emphasizing its advantage in terms of training stability. It is also worth noting that the repeated rollout mechanism in GRPO incurs higher API costs when interacting with a real search engine, further highlighting the practicality of our simulated search setting.

5.5 Reverse Curriculum Study

In this section, we analyze the effectiveness of the curriculum rollout strategy by comparing it with a reverse curriculum setup, in which the difficulty of training decreases over time by gradually improving the quality of retrieved documents. The results are presented in Table 6.

The results clearly indicate that the standard easy-to-hard curriculum consistently outperforms the reverse hard-to-easy variant across both models, supporting the effectiveness of curriculum learning in our framework. Starting with better search results allows the policy model to first learn how to invoke the search engine and understand the basic output format. As training progresses, the model is exposed to increasingly challenging scenarios, fostering stronger reasoning capabilities.

6 Conclusion

In this paper, we propose ZEROSEARCH, a novel RL framework that enhances the search capabilities of LLMs without interacting with real search engines. Through supervised fine-tuning, the LLM is transformed into a retrieval module capable of generating both relevant and noisy documents. A curriculum rollout mechanism is employed to progressively improve reasoning by exposing the model to increasingly challenging retrieval scenarios. Experimental results show that ZEROSEARCH outperforms real search-based models, generalizes well across both base and instruction-tuned LLMs of varying sizes, and supports a wide range of RL algorithms.

However, our approach has certain limitations. Deploying the simulated search LLM requires access to GPU servers. While more cost-effective than commercial API usage, this introduces additional infrastructure costs. We provide a detailed discussion of these costs in the appendix.

References

- [1] A. Asai, Z. Wu, Y. Wang, A. Sil, and H. Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*, 2023.
- [2] B. Bohnet, V. Q. Tran, P. Verga, R. Aharoni, D. Andor, L. B. Soares, J. Eisenstein, K. Ganchev, J. Herzig, K. Hui, et al. Attributed question answering: Evaluation and modeling for attributed large language models. *arXiv preprint arXiv:2212.08037*, 2022.
- [3] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [4] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [5] W. Feng, C. Hao, Y. Zhang, J. Song, and H. Wang. Airrag: Activating intrinsic reasoning for retrieval augmented generation via tree-based search. *arXiv preprint arXiv:2501.10053*, 2025.
- [6] L. Gao, Z. Dai, P. Pasupat, A. Chen, A. T. Chaganty, Y. Fan, V. Y. Zhao, N. Lao, H. Lee, D.-C. Juan, et al. Rarr: Researching and revising what language models say, using language models. *arXiv preprint arXiv:2210.08726*, 2022.
- [7] Y. Guo, L. Hou, R. Shao, P. G. Jin, V. Kumar, W. Weng, Y. Xie, and T.-Y. Liu. Deepseek-r1: Reinforcement learning for retrieval-augmented generation in large language models. *arXiv preprint arXiv:2503.01234*, 2025.
- [8] X. Ho, A.-K. D. Nguyen, S. Sugawara, and A. Aizawa. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*, 2020.
- [9] Y. Hou and et al. RL-based learning for reasoning and decision-making in large language models. In *ACL*, 2025.
- [10] J. Hu. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*, 2025.
- [11] S. Imani, L. Du, and H. Shrivastava. Mathprompter: Mathematical reasoning using large language models. *arXiv preprint arXiv:2303.05398*, 2023.
- [12] S. Jeong, J. Baek, S. Cho, S. J. Hwang, and J. C. Park. Adaptive-rag: Learning to adapt retrieval-augmented large language models through question complexity. *arXiv preprint arXiv:2403.14403*, 2024.
- [13] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023.
- [14] J. Jiang, J. Chen, J. Li, R. Ren, S. Wang, W. X. Zhao, Y. Song, and T. Zhang. Rag-star: Enhancing deliberative reasoning with retrieval augmented verification and refinement. *arXiv preprint arXiv:2412.12881*, 2024.
- [15] J. Jiang, Z. Chen, Y. Min, J. Chen, X. Cheng, J. Wang, Y. Tang, H. Sun, J. Deng, W. X. Zhao, et al. Technical report: Enhancing llm reasoning with reward-guided tree search. *arXiv preprint arXiv:2411.11694*, 2024.
- [16] Z. Jiang, F. F. Xu, L. Gao, Z. Sun, Q. Liu, J. Dwivedi-Yu, Y. Yang, J. Callan, and G. Neubig. Active retrieval augmented generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7969–7992, 2023.
- [17] B. Jin, H. Zeng, Z. Yue, D. Wang, H. Zamani, and J. Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025.

- [18] M. Joshi, E. Choi, D. S. Weld, and L. Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
- [19] R. Kumar and et al. Research: Autonomous retrieval decision-making in llms using reinforcement learning. In *ICLR*, 2025.
- [20] V. Kumar, L. Hou, Y. Guo, R. Shao, P. G. Jin, W. Weng, Y. Xie, and T.-Y. Liu. Self-correcting language models with reinforcement learning. *arXiv preprint arXiv:2409.06543*, 2024.
- [21] T. Kwiatkowski, J. Palomaki, O. Redfield, M. Collins, A. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
- [22] A. Lewkowycz, A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. Ramasesh, A. Slone, C. Anil, I. Schlag, T. Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857, 2022.
- [23] X. Li, G. Dong, J. Jin, Y. Zhang, Y. Zhou, Y. Zhu, P. Zhang, and Z. Dou. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366*, 2025.
- [24] X. Li, J. Jin, Y. Zhou, Y. Wu, Z. Li, Q. Ye, and Z. Dou. Retrollm: Empowering large language models to retrieve fine-grained evidence within generation. *arXiv preprint arXiv:2412.11919*, 2024.
- [25] X. Li, W. Xu, R. Zhao, F. Jiao, S. Joty, and L. Bing. Can we further elicit reasoning in llms? critic-guided planning with retrieval-augmentation for solving challenging tasks. *arXiv preprint arXiv:2410.01428*, 2024.
- [26] A. Mallen, A. Asai, V. Zhong, R. Das, H. Hajishirzi, and D. Khashabi. When not to trust language models: Investigating effectiveness and limitations of parametric and non-parametric memories. *arXiv preprint arXiv:2212.10511*, 7, 2022.
- [27] J. Menick, M. Trebacz, V. Mikulik, J. Aslanides, F. Song, M. Chadwick, M. Glaese, S. Young, L. Campbell-Gillingham, G. Irving, et al. Teaching language models to support answers with verified quotes. *arXiv preprint arXiv:2203.11147*, 2022.
- [28] O. Press, M. Zhang, S. Min, L. Schmidt, N. A. Smith, and M. Lewis. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*, 2022.
- [29] O. Ram, Y. Levine, I. Dalmedigos, D. Muhlgay, A. Shashua, K. Leyton-Brown, and Y. Shoham. In-context retrieval-augmented language models. *arXiv preprint arXiv:2302.00083*, 2023.
- [30] H. Rashkin, V. Nikolaev, M. Lamm, L. Aroyo, M. Collins, D. Das, S. Petrov, G. S. Tomar, I. Turc, and D. Reitter. Measuring attribution in natural language generation models. *arXiv preprint arXiv:2112.12870*, 2021.
- [31] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [32] Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [33] W. Shi, S. Min, M. Yasunaga, M. Seo, R. James, M. Lewis, L. Zettlemoyer, and W.-t. Yih. Replug: Retrieval-augmented black-box language models. *arXiv preprint arXiv:2301.12652*, 2023.
- [34] K. Shuster, S. Poff, M. Chen, D. Kiela, and J. Weston. Retrieval augmentation reduces hallucination in conversation. *arXiv preprint arXiv:2104.07567*, 2021.
- [35] H. Song, J. Jiang, Y. Min, J. Chen, Z. Chen, W. X. Zhao, L. Fang, and J.-R. Wen. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. *arXiv preprint arXiv:2503.05592*, 2025.

- [36] R. Taylor, M. Kardas, G. Cucurull, T. Scialom, A. Hartshorn, E. Saravia, A. Poulton, V. Kerkez, and R. Stojnic. Galactica: A large language model for science. *CoRR*, abs/2211.09085, 2022.
- [37] H. Trivedi, N. Balasubramanian, T. Khot, and A. Sabharwal. Musique: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554, 2022.
- [38] S. Xia, X. Li, Y. Liu, T. Wu, and P. Liu. Evaluating mathematical reasoning beyond accuracy. *arXiv preprint arXiv:2404.05692*, 2024.
- [39] R. Yamauchi, S. Sonoda, A. Sannai, and W. Kumagai. Lpml: llm-prompting markup language for mathematical reasoning. *arXiv preprint arXiv:2309.13078*, 2023.
- [40] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [41] Z. Yang, P. Qi, S. Zhang, Y. Bengio, W. W. Cohen, R. Salakhutdinov, and C. D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.
- [42] O. Yoran, T. Wolfson, B. Bogin, U. Katz, D. Deutch, and J. Berant. Answering questions by meta-reasoning over multiple chains of thought. *arXiv preprint arXiv:2304.13007*, 2023.
- [43] W. Yu, D. Iter, S. Wang, Y. Xu, M. Ju, S. Sanyal, C. Zhu, M. Zeng, and M. Jiang. Generate rather than retrieve: Large language models are strong context generators. *arXiv preprint arXiv:2209.10063*, 2022.
- [44] J. Zhang, Z. Li, K. Das, B. Malin, and S. Kumar. Sac3: Reliable hallucination detection in black-box language models via semantic-aware cross-check consistency: Reliable hallucination detection in black-box language models via semantic-aware cross-check consistency. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15445–15458, 2023.
- [45] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [46] Y. Zhao, H. Yin, B. Zeng, H. Wang, T. Shi, C. Lyu, L. Wang, W. Luo, and K. Zhang. Marco-o1: Towards open reasoning models for open-ended solutions. *arXiv preprint arXiv:2411.14405*, 2024.
- [47] Y. Zheng, D. Fu, X. Hu, X. Cai, L. Ye, P. Lu, and P. Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv preprint arXiv:2504.03160*, 2025.

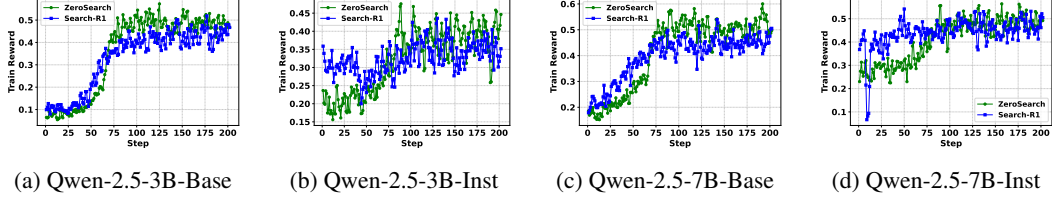


Figure 3: Reward curve comparison between ZEROSEARCH and Search-R1(using a real search engine).

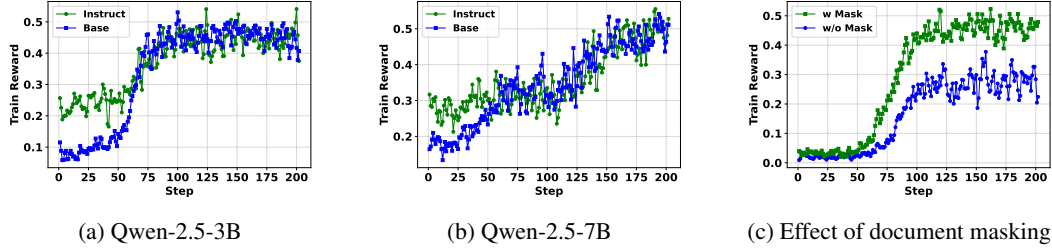


Figure 4: (a-b) We compare the reward curve between base and instruct models using Qwen-2.5-3B and Qwen-2.5-7B models. (c): We study the effects of document token loss masking using LLaMA-3.2-3B-base.

A Compare ZEROSEARCH with Real Search Engine

In this section, we present additional results comparing ZEROSEARCH with a real search engine using the Qwen-2.5-3B and Qwen-2.5-7B model series in Figure 3.

Across both model sizes, ZEROSEARCH consistently achieves a smoother reward curve compared to the real search engine. This is primarily because the quality of documents returned by the real search engine is uncontrollable during rollout. In the early stages of training, low-quality documents may prevent the policy model from developing a correct understanding of the task. In later stages, if the documents are too high-quality, the policy model may not be sufficiently challenged to continue improving its reasoning capability. In contrast, ZEROSEARCH enables dynamic control over document difficulty throughout training. This allows the policy model to first build a foundational understanding of the task and then gradually adapt to more complex scenarios.

B Compare Base and Instruct LLMs

In this section, we compare the training reward curves of base and instruction-tuned models using Qwen-2.5-3B and Qwen-2.5-7B. The results are presented in Figures 4a and 4b.

As shown, instruction-tuned models initially achieve higher rewards, owing to their stronger instruction-following capabilities, which allow them to invoke the search engine more effectively in the early stages of training. As training progresses, both base and instruction-tuned models demonstrate steady reward improvements. Notably, base models demonstrate greater reward improvements and reach performance levels comparable to their instruction-tuned counterparts. These results underscore the compatibility of ZEROSEARCH with both base and instruction-tuned models. Furthermore, they demonstrate that base models can effectively acquire search capabilities through reinforcement learning without the need for supervised fine-tuning as a warm-up.

C Effect of Document Token Loss Masking

During training, we apply loss masking to document tokens, as they are not generated by the policy model and may introduce noise. To assess the impact of the loss masking, we conduct ablation experiments using the LLaMA-3.2-3B model. The resulting reward curves are shown in Figure 4c.

Method	NQ	TriviaQA	PopQA	HotpotQA	2Wiki	Musique	Bamboogle	Avg.
w/ mask	40.28	61.87	49.80	30.57	35.05	14.75	14.52	35.26
w/o mask	34.70	54.58	44.51	23.08	30.41	10.71	13.01	30.14

Table 7: Effect of Document Token Loss Masking. We compare model performance with and without loss masking applied to document tokens.

Search Engine	Queries	Training Time	GPUs Used	API Cost	GPU Cost	Total Cost
SFT-3B	~64,000	~12 hours	1 × A100 GPUs	\$0.0	\$17.7	\$17.7
SFT-7B	~64,000	~12 hours	2 × A100 GPUs	\$0.0	\$35.4	\$35.4
SFT-14B	~64,000	~12 hours	4 × A100 GPUs	\$0.0	\$70.8	\$70.8
Google	~64,000	~12 hours	None	\$586.7	\$0.0	\$586.7

Table 8: Cost comparison between real search engine and our simulated search approach.

As illustrated, removing document token loss masking leads to a substantial drop in the reward, indicating that including such tokens in the loss computation causes training instability and degrades learning effectiveness. Table 7 further supports this observation, showing a significant decline in model performance without loss masking.

D Cost Analysis

Although our method eliminates the cost associated with commercial APIs, it requires the deployment of simulated search LLMs on GPU servers. In this section, we compare the cost of our method with that of commercial search API-based approaches. API costs are estimated based on the pricing of SerpAPI³, while GPU deployment costs are based on Amazon Web Services (AWS)⁴. In this paper, we perform reinforcement learning (RL) training with a batch size of 64, a rollout repetition of 5, and a total of 200 training steps. This setup results in approximately 12 hours of training when using Qwen-2.5-7B as the backbone, generating approximately 64,000 search requests.

As shown in Table 8, the cost of using Google Search is significantly higher than that of our approach. Furthermore, the API cost increases substantially with larger batch sizes or more training steps, making the method less scalable due to prohibitively high expenses. In contrast, ZEROSEARCH offers a more cost-efficient alternative. It is important to note that GPU utilization of the simulation server is inconsistent—high during the rollout and low during policy update—resulting in periods of underutilization and resource waste. A practical strategy to mitigate this issue is to share a single simulation server across multiple RL training tasks, which can further reduce deployment costs. In addition, we provide simulation LLMs of various parameter sizes, giving users flexible options to balance performance and resource usage based on their specific needs.

E Case Study

In this section, we present case studies from the interaction trajectories and the simulated documents to further illustrate the effectiveness of our proposed method.

Interaction Trajectory Study We first show several interaction trajectories in Table 9. From these examples, we make the following observations: First, the policy model consistently adheres to the expected output format, which is surprising given that the format requirements are only specified in the input template and are not explicitly reinforced through the reward design. Second, the model demonstrates the capability for multi-turn search behavior to arrive at the final answer. This is a crucial finding, as it confirms that our method effectively incentivizes the model’s search capabilities.

Simulated Document Study We further present examples of simulated document outputs generated by the simulation LLMs in Table 10. By comparing useful outputs with noisy ones, we observe a

³<https://serpapi.com/pricing>

⁴<https://aws.amazon.com/ec2/capacityblocks/pricing>

Question: Who is the spouse of the person who does the voice for Smokey the Bear?

Answer: Katharine Ross

ZEROSEARCH:

<think> I need to find the voice actor for Smokey the Bear and then find his spouse. **</think>**

<search> Smokey the Bear voice actor **</search>**

<information> Doc 1: Stephen Colbert, Al Roker and Jeff Foxworthy have joined the historic campaign, lending their voices to help expand on Smokey’s iconic catchphrase. Doc 2: TIL actor Sam Elliott – known for his full mustache, and deep, resonant voice – voices Smokey the Bear in the iconic character’s public Doc 3: Actor Sam Elliott, on sharing a birthday and voice with Smokey Bear With his bright, white mustache and earthy, western drawl, Sam Elliott Doc 4: Sam Elliot was a great Smokey Bear! Other actors who probably could have voiced him include Kevin Tighe and Randolph Mantooth whose character on Doc 5: Shares birthday with "Smokey Bear," and Sam Elliott has been the voice of "Smokey Bear" since 2008. Inducted into the International Mustache Hall of Fame in **</information>**

<think> The voice actor for Smokey the Bear is Sam Elliott. Now I need to find his spouse. **</think>**

<search> Sam Elliott spouse **</search>**

<information> Doc 1: Sam Elliot and Katharine Ross have been married for 35 years. Here’s how they keep the spark alive even when working together on a film. Doc 2: Meet Sam Elliott’s And Katharine Ross’ 40-Year Love Story#hollywood · Comments1.8K. Doc 3: Sam Elliott and Katharine Ross met in 1978 and were married in 1984. They met while filming the horror movie, The Legacy. It was Ross’s fifth Doc 4: Sam Elliott and his wife Katharine Ross have been married for nearly four decades. Their love story dates back to the late 1960s. Doc 5: 1978 Pete (Sam Elliott) and Maggie (Katherine Ross) are lovers to whom strange things begin to happen, in Universal’s horror thriller ‘The Legacy’ **</information>**

<think> Sam Elliott’s spouse is Katharine Ross. **</think>**

<answer> Katharine Ross **</answer>**

Question: What is the official color of the third oldest surviving university?

Answer: Cambridge Blue

ZEROSEARCH:

<think> I need to find the third oldest surviving university and its official color. **</think>**

<search> third oldest surviving university **</search>**

<information> Doc 1: Founded in 1209, Cambridge is the world’s third oldest surviving university and one of its most prestigious, currently ranked second best in the Doc 2: The oldest university in the English-speaking world, Oxford is regularly recognised as one of the premier academic institutions on the planet. Doc 3: The oldest education institution of any kind is likely The King’s School, Canterbury, which has been in continuous operation since 597. Doc 4: Founded in 1088, the University of Bologna is the oldest surviving university in Europe and one of the most prestigious universities in Italy. Doc 5: The University of Oxford is the third oldest university in the world, established in the twelfth century in Oxford, England. With an esteemed **</information>**

<think> The third oldest surviving university is Cambridge University. Now I need to find its official color. **</think>**

<search> official color of cambridge university **</search>**

<information> Doc 1: The logo only appears in the four colour variants shown on this page. Black, Red Pantone 032, Yellow Pantone 109 and white coat of arms with Doc 2: Official University of Cambridge Large Crest Print T-Shirt – Green. £19.99. T-shirt with a one-colour print crest and Cambridge University design. Part of the Doc 3: “The colour is called Cambridge Blue.” “It doesn’t matter what it’s called. It LOOKS green so – it IS green,” she added. I was perplexed and Doc 4: Premium quality unisex relaxed-fit t-shirt with detailed neck-tape and label, embroidered with the University of Cambridge crest. Discontinued colour Doc 5: Cambridge Blue is a soft, muted shade of blue-green that exudes a sense of calm and sophistication. This distinctive color is named after the University of **</information>**

<think> The official color of Cambridge University is Cambridge Blue. **</think>**

<answer> Cambridge Blue **</answer>**

Table 9: Interaction trajectories of ZEROSEARCH. The results demonstrate that the policy model learns to output the correct format and effectively conducts multi-step search to reach the final answer.

clear distinction in quality. Notably, the useful outputs consistently contain the correct answer, while the noisy outputs fail to do so. Precise control over document quality enables the implementation of a curriculum-based rollout mechanism and contributes to stabilizing the reinforcement learning training process.

Question: When was the director of film Reykjavík-Rotterdam born?

Answer: 30 June 1963

Query: Óskar Jónasson birth date

Useful Output:

Doc 1: Óskar Jónasson is a television and film director from Iceland. His feature debut Reykjavík-Rotterdam won the Golden Bear at the Berlin International Film Festival

Doc 2: Director of Reykjavík-Rotterdam and Silence of My Father Óskar Jónasson was born on **30 June 1963** in Reykjavík, Iceland. He is a graduate of the Iceland Academy of

Doc 3: A documentary about Iceland's aborted attempt to build a nuclear power plant in the 1970s, as a result of which the nation became largely

Doc 4: Director of Reykjavík-Rotterdam and Silence of My Father Óskar Jónasson was born on **30 June 1963** in Reykjavík, Iceland. He is a graduate of the Iceland Academy

Doc 5: Icelandic filmmaker Oscar Jonasson talks about his new film, the Golden Bear-winning "Reykjavik-Rotterdam," at the Berlinale this year.

Noisy Output:

Doc 1: The first single, "Láta ég allt" was released on March 18th 2022. It was written by him along with Óskar Jónasson and is accompanied by music video by the multi-

Doc 2: The children of the first generation include the likes of Ásgeir Trausti (born 1984) and Amiina (formed 2005), who also introduced the world

Doc 3: Óskar Jónasson · 2022 · Iceland · Tracks: 13 · Genres: Pop · Languages: Icelandic.

Doc 4: Óskar Jónasson. · Location. Reykjavík, Iceland. · Industry. Music. · Looking for. Jobs. · Experience. 7 years. · Education. 4 years.

Doc 5: Jónasson was born in Iceland, and his family moved to Sweden when he was a child. He graduated from Stockholm University with a degree in

Question: What is the official color of the third oldest surviving university?

Answer: Cambridge Blue

Query: third oldest surviving university

Useful Output:

Doc 1: The University of Oxford is the oldest university in the English-speaking world, with origins that can be traced back to 1096.

Doc 2: **The University of Cambridge** is the third oldest surviving university in the world, after the University of Bologna in Italy and the University of Paris in

Doc 3: **University of Cambridge**. The University of Cambridge is the third oldest surviving university in the world, after the University of Bologna in Italy and the

Doc 4: The University of Oxford is the oldest university in the English-speaking world, with origins that can be traced back to 1096.

Doc 5: **The University of Cambridge** is the third oldest surviving university in the world, after the University of Bologna in Italy and the University of Paris in

Noisy Output:

Doc 1: The oldest continuously operating university in the world is the University of Bologna, Italy. Founded in 1088, it is the oldest of the

Doc 2: Founded in 1088, the University of Bologna is the oldest in continuous operation. It is also among the most respected universities in the world.

Doc 3: The University of Bologna is the oldest in the world, founded in 1088, and the University of Oxford and the University of Cambridge are both

Doc 4: The University of Bologna is the oldest surviving university in the world, founded in 1088. It is also one of the most prestigious, having educated some of the

Doc 5: Founded in 1088, the University of Bologna is the oldest still in operation. The University of Paris, founded around 1150, is the second oldest.

Table 10: Outputs from the 14B simulation LLMs. The correct answers are highlighted in red. We can find that the quality of the useful output is much better than that of the noisy output.