# WEBEXPLORER: EXPLORE AND EVOLVE FOR TRAIN-ING LONG-HORIZON WEB AGENTS

# **Anonymous authors**

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

034

039

045

046 047 048

051

052

Paper under double-blind review

#### **ABSTRACT**

The paradigm of Large Language Models (LLMs) has increasingly shifted toward agentic applications, where web browsing capabilities are fundamental for retrieving information from diverse online sources. However, existing open-source web agents either demonstrate limited information-seeking abilities on complex tasks or lack transparent implementations. In this work, we identify that the key challenge lies in the scarcity of challenging data for information seeking. To address this limitation, we introduce WEBEXPLORER: a systematic data generation approach using model-based exploration and iterative, long-to-short query evolution. This method creates challenging query-answer pairs that require multistep reasoning and complex web navigation. By leveraging our curated highquality dataset, we successfully develop advanced web agent WEBEXPLORER-8B through supervised fine-tuning followed by reinforcement learning. Our model supports 128K context length and up to 100 tool calling turns, enabling longhorizon problem solving. Across diverse information-seeking benchmarks, WE-BEXPLORER-8B achieves the state-of-the-art performance at its scale. Notably, as an 8B-sized model, WEBEXPLORER-8B is able to effectively search over an average of 16 turns after RL training, achieving higher accuracy than WebSailor-72B on BrowseComp-en/zh and attaining the best performance among models up to 100B parameters on WebWalkerQA and FRAMES. Beyond these informationseeking tasks, our model also achieves strong generalization on the HLE benchmark even though it is only trained on knowledge-intensive QA data. These results highlight our approach as a practical path toward long-horizon web agents.

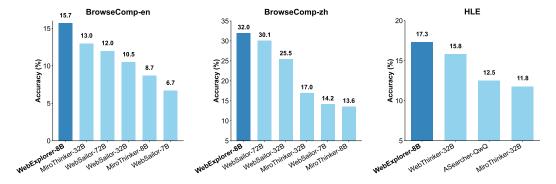


Figure 1: Performance comparison on BrowseComp-en, BrowseComp-zh and HLE benchmarks across different models.

# 1 Introduction

The rapid advancement of Large Language Models (LLMs) has spurred a paradigm shift toward agentic intelligence, where models are equipped with diverse tools to tackle increasingly complex problems. Web browsing agents (OpenAI, 2025; Perplexity AI, 2025) represent a critical component of this shift, enabling autonomous information retrieval from the vast landscape of online sources. Despite significant progress, existing open-source models such as WebThinker (Li et al., 2025b)

056

057

058

059

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

077

079

080

081

082

083

084

085

086

087

880

089

090

091 092

093

094

095

096

098

100

101

102

103

104

105

106

107

and WebSailor (Li et al., 2025a) exhibit limited performance on challenging search-based benchmarks like BrowseComp-en/zh (Wei et al., 2025; Zhou et al., 2025), BrowseComp-Plus (Chen et al., 2025) and Xbench (Xbench-Team, 2025), while stronger commercial models remain opaque in their training details (Zhipu AI, 2025; OpenAI, 2025).

We hypothesize that the fundamental challenge in developing capable web agents lies in the quality of training data. Current benchmarks have evolved to include queries of extraordinary difficulty – For example, over half of BrowseComp-en queries fail to be solved by human annotators. While such challenging queries are rare in typical use cases, constructing high-quality, difficult query-answer pairs is essential for developing agents that can achieve super-human performance on information-seeking tasks. Previous approaches to web navigation data construction have primarily employed two methodologies: Graph-based methods (Li et al., 2025a; Tao et al., 2025) construct explicit graphs through random walks where each node represents a website or entity and each edge represents a relationship, then utilize subgraphs to formulate QA pairs. Evolution-based approaches (Wu et al., 2025a; Gao et al., 2025) iteratively modify existing queries to increase their complexity. However, these approaches face limitations. Previous graph construction introduces complexity in node expansion and selection heuristics. Meanwhile, evolution-based methods usually increase the complexity by injecting new information to replace original content at each evolution step, potentially resulting in unnatural query formulations and limited flexibility in synthesis. Our framework WEBEXPLORER addresses these limitations through a novel approach that leverages model-based exploration to construct the information space and iterative long-to-short query evolution to reduce salient information. This enables more natural and flexible synthesis of challenging web navigation tasks.

Rather than constructing explicit graphs through rule-based methods, which is a complex process requiring decisions about expansion strategies and node selection, we adopt a simpler model-based approach to explore the information space purely through prompting. Starting from a seed entity, we leverage powerful LLMs to simulate the graph-building process internally through iterative search and browsing actions. This approach enables flexible, dynamic exploration of information spaces related to the seed entity without the overhead of explicit graph construction. The model then utilizes this explored information space to construct initial query-answer pairs.

In our preliminary experiments, however, we observed that these initially constructed QA pairs, while requiring multiple websites to solve, still remained relatively straightforward evidenced by the high success rates achieved by strong proprietary models. To address this limitation, we introduce a second-stage evolution process. Specifically, unlike previous work where evolution involves injecting new information (Wu et al., 2025a; Gao et al., 2025), we prompt models to systematically increase query difficulty by *removing* explicit clues and introducing strategic obfuscation. In contrast to the initial QA pairs with clear search entry points, the evolved QA pairs require longer solution processes with more exploratory search attempts. This systematic evolution generates challenging queries that lead to significantly lower success rates and require more reasoning steps from proprietary models. Through this process, we construct our WEBEXPLORER-QA data.

Our training utilizes a typical approach combining supervised fine-tuning for cold-start initialization, followed by reinforcement learning using the GRPO algorithm (Shao et al., 2024). Our RL training scales to 128K context length and 100 maximum number of tool calling turns, where we observe consistent increases in both the number of tool calls and benchmark performance. Based on Qwen3-8B (Yang et al., 2025), our WEBEXPLORER-8B achieves state-of-the-art performance at its scale on multiple information-seeking benchmarks, including BrowseComp-en/zh (Wei et al., 2025; Zhou et al., 2025), GAIA (Mialon et al., 2024), WebWalkerQA (Wu et al., 2025b), Frames (Krishna et al., 2024), and XBench-DeepSearch (Xbench-Team, 2025). Notably, our model achieves 15.7% on BrowseComp-en and 32.0% on BrowseComp-zh, significantly outperforming the previous leading WebSailor 72B model despite a much smaller size. It also achieves 62.7% on WebWalkerQA and 75.7% on FRAMES, establishing the best performance among models up to 100B parameters. Beyond superior performance on information-seeking tasks, WEBEXPLORER-8B demonstrates remarkable generalization to the academic benchmark HLE (Phan et al., 2025), scoring 17.3% and outperforming previous 32B models like WebThinker-32B (Li et al., 2025b). This validates the strong generalization capability of our approach beyond pure information-seeking tasks. Overall, the superior performance of WEBEXPLORER-8B across diverse benchmarks strongly validates the quality of our synthesized data and establishes a practical pathway toward building capable longhorizon web agents.

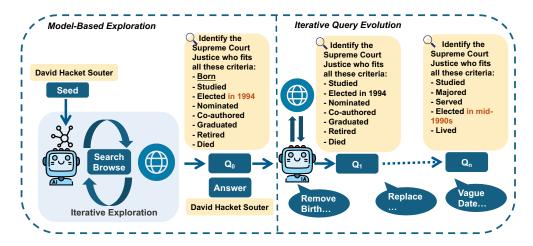


Figure 2: Model-Based Exploration and Iterative Query Evolution Framework. Starting from a seed entity (e.g., David Hackett Souter), the framework employs iterative search and browsing actions to construct the information space related to the seed entity. Initial queries  $(Q_0)$  and Answers are generated based on this explored information space. Through iterative evolution, salient information is systematically obfuscated (e.g., "Remove Birth...", "Replace ..." or "Vague Date...") to produce more challenging queries  $(Q_1$  to  $Q_n)$ . This process ensures the resulting queries require longer reasoning steps and explorations.

# 2 WebExplorer: Synthesizing Challenging QA Pairs

Contemporary information-seeking benchmarks present great challenges, exemplified by BrowseComp-en (Wei et al., 2025) where more than half of the queries remain unsolvable by human annotators despite two hours of dedicated effort. These challenging benchmarks rely predominantly on manual curation (Wei et al., 2025; Zhou et al., 2025; Mialon et al., 2024), resulting in prohibitively high annotation costs and severely limited dataset sizes that preclude their use in large-scale training. Concurrently, existing open-source training data exhibit some limitations: they either lack sufficient scale to enable large-scale model training or present overly simplistic queries that fail to match the difficulty of evaluation benchmarks (Wu et al., 2025a; Li et al., 2025a; Tao et al., 2025; Pahuja et al., 2025). Consequently, the autonomous synthesis of large-scale, challenging search-oriented query-answer pairs emerges as a pivotal research challenge in developing advanced web agents.

Several works have explored synthesizing challenging information-seeking QA pairs through two predominant methodological paradigms. Graph-based approaches explicitly construct knowledge graphs where each node represents a website or entity, then generate QA pairs by selecting subgraphs from the complete graph structure. Typically, larger subgraphs correspond to more challenging queries (Li et al., 2025a; Tao et al., 2025). Evolution-based methods iteratively modify existing queries by injecting new information to systematically increase their complexity (Wu et al., 2025a; Gao et al., 2025), typically evolving queries to become progressively longer. However, the previous rule-based graph construction introduces complexity of node expansion and selection heuristics. Instead, our method Webexplorer utilizes existing LLMs to navigate the web and construct the information space autonomously. Additionally, traditional evolution approaches that inject new information at every step can become unnaturally complex. Rather than following short-to-long evolution, our iterative query evolution operates in a *long-to-short* manner by strategically removing salient information, creating natural obfuscation. As illustrated in Figure 2, through these two stages, Webexplorer proposes a simple and flexible prompting-based framework for challenging QA pair generation.

#### 2.1 PRELIMINARIES

Our agent framework builds upon the ReAct paradigm (Yao et al., 2023), where agents execute iterative cycles of reasoning and acting. When presented with a query, the agent conducts multiple rounds of Thought-Action-Observation sequences. During each round, the language model

formulates a Thought  $\tau_t$  based on the current context, executes a structured Action  $\alpha_t$ , and receives corresponding environmental feedback as Observation  $o_t$ . The final round contains only Thought  $\tau_t$  without further actions and observations.

Our action space encompasses two primary tools:  $\alpha_t \in \{\texttt{search}, \texttt{browse}\}$ . The search action enables querying search engines to retrieve relevant information, while the <code>browse</code> action facilitates direct webpage access for content extraction. We formalize each action as  $\alpha_t = (\alpha_t^m, \alpha_t^p)$ , where  $\alpha_t^m$  specifies the action type and  $\alpha_t^p$  contains the necessary parameters. For <code>search</code> operations,  $\alpha_t^p$  encompasses a list of query strings, whereas <code>browse</code> operations require target URLs and extraction objectives within  $\alpha_t^p$ . More details about these tools can be found in §4.1.

The complete execution trajectory spanning T iterations is represented as:

$$H_T = (\tau_0, \alpha_0, o_0, \tau_1, \alpha_1, o_1, \dots, \tau_T)$$
(1)

At each timestep t, the agent produces thought  $\tau_t$  and selects action  $\alpha_t$  according to policy  $\pi(\tau_t, \alpha_t | H_{t-1})$ , conditioning on the accumulated trajectory history  $H_{t-1}$ . We provide a detailed trajectory example in Appendix A.

# 2.2 MODEL-BASED EXPLORATION

Traditional graph-based methods (Li et al., 2025a; Tao et al., 2025) require heuristic rules of explicit graph construction, typically initiated from a root URL or seed entity, followed by systematic expansion through predefined traversal strategies. The resulting graph contains structured information where usually nodes represent entities and edges form the relationships. Such approaches involve iteratively identifying related entities, extracting their features, and expanding the graph until reaching predetermined size constraints. This process introduces some complexity, requiring careful, heuristic design of expansion strategies and node selection heuristics.

In contrast, we propose WEBEXPLORER, which employs a different model-based exploration approach by leveraging powerful LLMs to construct the information space autonomously. Our method operates purely through prompting: we provide a seed entity as the initial search entry along with three example QA pairs, then instruct the model to conduct iterative search and browsing actions to explore the information space before synthesizing QA pairs. Specifically, given an entity as a seed, the models conduct iterative search and browsing actions to construct a comprehensive information space encompassing entity-related content. This approach simulates the graph-building process internally, enabling flexible and dynamic exploration without the complexity of graph expansion strategy design and explicit graph maintenance. Subsequently, utilizing the explored information space, models generate QA pairs that necessitate reasoning across multiple websites to reach the correct solution. A detailed example demonstrating how the model iteratively explores and then forms a query-answer pair using multiple sources of searched information is provided in Appendix C.

Formally, starting from a seed entity  $e_0$ , the exploration process conducts multiple rounds of actions and observations:

$$H_T = (e_0, \tau_0, \alpha_0, o_0, \tau_1, \alpha_1, o_1, \dots, \tau_T)$$
(2)

where H represents the complete information space encompassing all explored content, and  $\tau_T$  contains the synthesized QA pair. The model autonomously determines when to stop searching to formulate a challenging QA pair, eliminating the need for predefined stopping criteria and allowing the model to explore flexibly.

# 2.3 Iterative Query Evolution

# Example of BrowseComp-en

**Query:** Please identify the fictional character who occasionally breaks the fourth wall with the audience, has a backstory involving help from selfless ascetics, is known for his humor, and had a TV show that aired between the 1960s and 1980s with fewer than 50 episodes. *Answer: Plastic Man* 

While the initial QA pairs constructed from model-explored information spaces successfully incorporate content from multiple websites, we observe that proprietary models can still solve them with relatively high accuracy. For instance, as shown in Table 1, Claude-4-Sonnet achieves 86.6% accuracy on the initial QA pairs compared to 12.2% on BrowseComp-en and 68.3% on GAIA. This indicates that despite requiring multi-website reasoning, these initial QA pairs remain insufficiently challenging. Through systematic case analysis, we identify that the initial queries contain excessive explicit information and salient clues that substantially reduce their difficulty. For example, the initial query-answer example shown below demonstrates several salient clues such as "the official attendance set a record" and "this player died at the age of 44" that serve as clear and strong indicators. Such specific information—including dates, locations, and proper names—often provides direct entry points that enable straightforward solution trajectories without requiring exploratory detours or consideration of alternative reasoning paths.

In contrast, examination of challenging information-seeking benchmarks like BrowseComp-en reveals a critical distinguishing characteristic: these queries deliberately avoid providing clear, specific clues, instead employing vague descriptions. As the example above shows, the BrowseComp-en query "TV show that aired between the 1960s and 1980s with fewer than 50 episodes" demonstrates this obfuscation. Therefore, these challenging queries usually do not contain clear search entry points, which necessitate extensive exploration and multiple reasoning attempts before reaching the correct answer.

Drawing inspiration from BrowseComp's design principles, we implement an iterative query evolution process to systematically increase QA pair difficulty. Unlike previous evolution methods that follow a short-to-long approach (Wu et al., 2025a; Gao et al., 2025), our evolution primarily reduces excessive information from initial QA pairs to increase difficulty. Given the complete initial QA construction information, we explicitly prompt models to refine queries through three strategic directions: (1) removing salient information, (2) introducing strategic obfuscation to specific details such as dates, locations, and proper names, and (3) identifying alternative descriptions to replace original explicit references. Through multiple iterative evolution cycles, we generate substantially more challenging QA pairs that require extensive exploration and demonstrate significantly reduced success rates. As demonstrated in the example below, the evolved query eliminates some salient descriptions in purple color like *this player died at the age of 44* and obscures concrete entities by replacing explicit references such as *Manchester United* with vague descriptors like *First Division giant*.

# **Example of Evolved Query**

**Initial Query:** A football match took place in a stadium where the official attendance set a record that still stands today for FIFA World Cup matches. The referee of this match was the oldest person to ever officiate a World Cup final, and exactly 26 years after this match, he was the chairman of a club that defeated Manchester United in an FA Cup final. The player who scored the winning goal in that FA Cup final was born in an area that became part of its current city in 1920, and this player died at the age of 44. In what minute of the FA Cup final was the winning goal scored? *Answer:* 83rd minute

**Evolved Query:** In the unique FIFA World Cup tournament format that concluded without a knockout final, a match official later guided a Second Division club to victory over a First Division giant in the monarch's final attendance at such an occasion. The match-winner had been rejected by the club he supported as a child, hailing from a district that joined a centuries-old Royal Naval stronghold two decades into the 20th century. In which minute did this decisive strike occur? **Answer:** 83rd minute

Formally, starting with the initial query  $Q_0$  extracted from the previous information space  $H^{(0)}$ , we iteratively evolve queries through multiple rounds. For  $k=0,\ldots,K-1$ , the (k+1)-th information space  $H^{(k+1)}$  is obtained by appending the k-th round's augmentation to the entire history:

$$H^{(k+1)} = (H^k, \tau_1^{(k)}, \alpha_1^{(k)}, o_1^{(k)}, \dots, \tau_{m_k}^{(k)}), \quad k = 0, \dots, K - 1$$
(3)

where each evolved query  $Q_{k+1}$  is contained in  $\tau_{m_k}^{(k)}$  and maintains the same answer A. The sequence  $(\alpha_1^{(k)}, o_1^{(k)}, \dots, \alpha_{m_k-1}^{(k)}, o_{m_k-1}^{(k)})$  represents the multiple action-observation cycles conducted during the k-th evolution step to further augment the information space.

Dataset	Initial QA	Evolved QA	WebDancer	SailorFog	WebShaper	ASearcher
Average Turns	7.9	9.9	5.4	8.2	8.4	6.5
Accuracy (%)	86.6	67.1	62.0	35.0	67.4	62.0

Table 1: Comparison of average tool calling turns and accuracy (%) of Claude-4-Sonnet across web navigation datasets, reporting our Initial QA, Evolved QA, and other QA datasets (WebDancer, SailorFog, WebShaper, ASearcher).

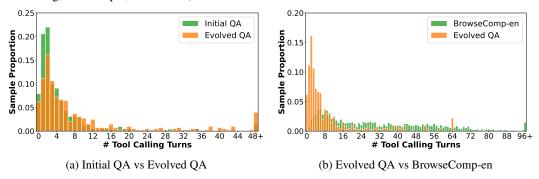


Figure 3: Tool calling turns distribution comparisons using OpenAI-o3: Initial QA vs Evolved QA (left) and Evolved QA vs BrowseComp-en (right).

# 2.4 THE WEBEXPLORER-QA DATASET

To ensure diversity in our synthesized QA pairs, we collect seed entities from Wikipedia and incorporate three QA pairs from BrowseComp-en as exemplars in our synthesis prompt. After generating initial QA pairs using our model-based exploration, we apply our iterative query evolution methodology for 5 iterations to synthesize the WEBEXPLORER-QA dataset, with detailed prompts description for the synthesis shown in Appendix B. These evolved final QA pairs, approximately 40K in total, are subsequently used for supervised fine-tuning and can be used directly for reinforcement learning.

To validate the quality of WEBEXPLORER-QA and demonstrate the effectiveness of our iterative evolution process, we conduct a comprehensive comparative analysis across multiple web navigation datasets. Our evaluation encompasses our initial QA pairs, the final evolved QA pairs, and established datasets from prior work, including WebDancer (Wu et al., 2025a), SailorFog (Li et al., 2025a), WebShaper (Tao et al., 2025), and ASearcher (Gao et al., 2025). Using Claude-4-Sonnet as our evaluation model, we assess both accuracy and the average number of tool calling turns required to solve each QA pair. Additionally, we employ a more powerful model, OpenAI-o3, to analyze the tool calling turns distribution between our Initial QA and Evolved QA, as well as between Evolved QA and the challenging benchmark BrowseComp-en.

The results in Table 1 demonstrate the critical importance of our iterative evolution process. Evolution significantly increases the complexity of the queries, with accuracy dropping from 86.6% to 67.1% and average solution turns increasing from 7.9 to 9.9, indicating successful creation of complex multi-step reasoning tasks. Furthermore, our evolved WEBEXPLORER-QA achieves the highest average turn count compared to existing datasets, demonstrating superior complexity. Figure 3 (left) also shows that easy QA pairs solvable within 4 turns decrease significantly after evolution. While Figure 3 (right) shows a gap remains between our Evolved QA and BrowseComp-en in tool calling turns, BrowseComp-en presents excessive difficulty with below 20% accuracy on most open-source models. Therefore, completely mirroring BrowseComp-en's difficulty level is unnecessary for training current open-source models.

### 3 COLD START AND REINFORCEMENT LEARNING

Our training methodology adopts the established post-training two-phase paradigm: supervised fine-tuning for cold start initialization, followed by reinforcement learning for advanced capability development (Guo et al., 2025; Zhipu AI, 2025). The initial supervised fine-tuning phase enables models to acquire proper invocation of search and browse functions while developing foundational

long-horizon search capabilities. Subsequently, reinforcement learning further enhances reasoning abilities, extending model performance to longer contexts and increased maximum turn limits, ultimately achieving more advanced long-horizon problem-solving behaviors.

### 3.1 SUPERVISED FINE-TUNING FOR COLD START

After synthesizing the challenging QA pairs as described in §2.4, we leverage the commercial model to collect high-quality trajectories for supervised fine-tuning. We employ rejection sampling during data collection, ensuring that our fine-tuning process exclusively utilizes correct trajectories. We adopt the ReAct framework (Yao et al., 2023) as our foundational format, incorporating search and browse as the two primary tools for actions  $\alpha$ , enclosed by <tool\_call> and </tool\_call> tags. The framework includes explicit reasoning thoughts  $\tau$  denoted by <think> tags, enabling transparent cognitive processes, and tool responses as observations  $\sigma$  marked by <tool\_response> tags. Our collected trajectories encompass multiple rounds of action  $\sigma$ , thought  $\tau$ , and observation  $\sigma$  sequences, with an example demonstrated in Appendix A.

#### 3.2 Reinforcement Learning

Following the cold-start phase, which endows the model with fundamental search and browse action capabilities along with long-horizon reasoning ability, we conduct further reinforcement learning training to enhance reasoning performance and optimize decision-making strategies using GRPO algorithm (Shao et al., 2024). Notably, in the RL phase, we can directly utilize the synthesized QA pairs without requiring solving trajectories. For reward design, we implement a composite reward function that balances structural correctness with answer accuracy:

$$R = 0.2 \cdot R_{\text{format}} + R_{\text{correct}} \tag{4}$$

Format rewards  $R_{\rm format}$  evaluate the correctness of response formatting, primarily assessing whether tool calls and thought structures adhere to the specified format requirements. For accuracy rewards  $R_{\rm correct}$ , we leverage the DeepSeek-V3 model (Liu et al., 2024) as an automated judge to evaluate whether the final responses are correct given the ground truth answers. This automated evaluation approach enables scalable assessment while maintaining high reliability.

To accommodate the model's expanding reasoning capabilities during training, we implement a progressive context length expansion strategy. We begin with a maximum length of 64K tokens and a tool calling turn limit of 50. As the model generates increasingly complex trajectories, we gradually increase the maximum length to 96K tokens with a 75-turn limit, and finally to 128K tokens with a 100-turn limit. This progressive expansion allows the model to develop more sophisticated long-horizon reasoning patterns throughout the training process.

#### 4 EXPERIMENTS

# 4.1 EXPERIMENTAL SETUPS

**Benchmarks** We evaluate on several information-seeking benchmarks including BrowseCompen (Wei et al., 2025) and BrowseComp-zh (Zhou et al., 2025), GAIA (Mialon et al., 2024) (103-sample text-only validation subset), WebWalkerQA (Wu et al., 2025b), FRAMES (Krishna et al., 2024), and XBench-DeepSearch (Xbench-Team, 2025). We also evaluate on the frontier academic benchmark HLE (Phan et al., 2025) to verify generalization capabilities beyond Wikipedia-alike knowledge QA. We report Avg@4 scores using LLM-as-Judge evaluation with DeepSeek-V3 (Liu et al., 2024) following previous work (Li et al., 2025a; Tao et al., 2025).

**Models** We compare our approach against both proprietary and open-source agents. For proprietary models, we primarily benchmark against OpenAI-o3 and Claude-4-Sonnet (Anthropic, 2025). Among open-source agents, we compare against advanced models including, GLM-4.5 (Zhipu AI, 2025), DeepSeek-V3.1 (Liu et al., 2024), Kimi-K2 (Team et al., 2025), WebShaper (Tao et al., 2025), WebSailor (Li et al., 2025a), ASearcher (Gao et al., 2025), WebThinker (Li et al., 2025b) and MiroThinker (MiroMind Team, 2025).

Model	BC-en	BC-zh	GAIA	WebWalkerQA	FRAMES	Xbench-DS	HLE
OpenAI-o3 <sup>†</sup>	50.9	58.1	70.5 <sup>†</sup>	71.7	84.0	66.7	20.2
Claude-4-Sonnet <sup>†</sup>	12.2	29.1	$68.3^{\dagger}$	61.7	80.7	64.6	20.3
GLM-4.5	26.4	37.5	$66.0^{\dagger}$	$65.6^{\dagger}$	$78.9^{\dagger}$	$70.0^{\dagger}$	21.2 <sup>†</sup>
DeepSeek-V3.1	30.0	49.2	$63.1^{\dagger}$	$61.2^{\dagger}$	83.7	71.2	29.8
Kimi-K2 <sup>†</sup>	14.1	28.8	57.7	63.0	72.0	50.0	18.1
WebShaper-72B	-	-	60.0	52.2	-	-	_
WebShaper-32B (QwQ)	-	-	53.3	49.7	-	-	-
WebShaper-32B	-	-	52.4	51.4	-	-	-
WebSailor-72B	12.0	30.1	55.4	-	-	55.0	-
WebSailor-32B	10.5	25.5	53.2	-	-	53.3	-
WebSailor-7B	6.7	14.2	33.0	-	-	34.3	-
ASearcher-Web-QwQ	5.2	15.6	52.8	34.3	70.9	42.1	12.5
WebThinker-32B	2.8	-	48.5	46.5	-	-	15.8
MiroThinker-32B-DPO-v0.1	13.0	17.0	57.3	49.3	71.7	-	11.8
MiroThinker-8B-DPO-v0.1	8.7	13.6	46.6	45.7	64.4	-	
WebExplorer-8B (SFT)	7.9	21.3	43.7	59.8	72.6	47.5	16.0
WebExplorer-8B (RL)	<u>15.7</u>	<u>32.0</u>	<u>50.0</u>	<u>62.7</u>	<u>75.7</u>	<u>53.7</u>	<u>17.3</u>

Table 2: Accuracy (%) of web agents on information-seeking benchmarks. BC-en and BC-zh denote BrowseComp-en and BrowseComp-zh respectively. XBench-DS refers to XBench-DeepSearch. **Bold** indicates the best performance among open-source models < 100B, while <u>underlined</u> values represent the best performance among models < 10B parameters. All scores of WEBEXPLORER-8B are computed as Avg@4 using LLM-as-Judge. Entries marked with a dagger (†) were reproduced by us under our scaffold: on model name = entire row; on a number = that entry only.

**Scaffold Details** Following previous work (Li et al., 2025a; Tao et al., 2025), our agent scaffold uses two tools: search returns top-10 Google results and browse retrieves URL content via Jina (Jina.ai, 2025) and answers queries using Gemini 2.5 Flash (Comanici et al., 2025). Further details about these two tools can be found in Appendix D. Using this unified scaffold, we evaluate OpenAI-o3, Claude-4-Sonnet, and Kimi-K2 on all benchmarks, and evaluate GLM-4.5 and DeepSeek-V3.1 on a subset (marked with † where the original reports lack results), ensuring a fair cross-model comparison.

**Training Details** We use the Qwen3-8B model (Yang et al., 2025) to conduct SFT and RL training, resulting in the WEBEXPLORER-8B model. For supervised fine-tuning, we utilize approximately 13K training samples with a training batch size of 32 and a learning rate of 1e-5 for 4 epochs. For reinforcement learning, we train on approximately 12K samples using the GRPO algorithm, where each group consists of 8 rollouts with a batch size of 64 and a learning rate of 1e-6.

### 4.2 EXPERIMENTAL RESULTS

WEBEXPLORER-8B Establishes State-of-the-Art Performance at Its Scale As shown in Table 2, our WEBEXPLORER-8B demonstrates exceptional performance across diverse information-seeking benchmarks. Following supervised fine-tuning, WEBEXPLORER-8B (SFT) already exhibits competitive results, surpassing WebSailor-7B with scores of 7.9% on BrowseComp-en, 72.6% on FRAMES, and 47.5% on XBench-DS. The subsequent reinforcement learning phase yields substantial improvements, with our final model achieving 15.7% on BrowseComp-en and 32.0% on BrowseComp-zh. Notably, WEBEXPLORER-8B (RL) surpasses the previous best sub-10B model, MiroThinker-8B-DPO-v0.1, across all benchmarks, establishing new state-of-the-art performance at its parameter scale.

**WEBEXPLORER-8B Outperforms Much Larger Models** Despite having only 8B parameters, WEBEXPLORER-8B consistently outperforms much larger open-source models, demonstrating remarkable parameter efficiency. On BrowseComp-en, BrowseComp-zh, FRAMES, and WebWalk-

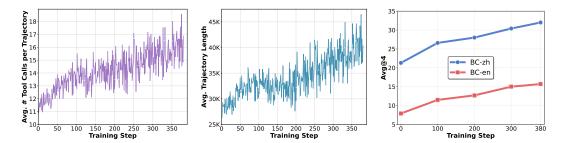


Figure 4: Training progress during RL. **Left**: Average # tool calls per trajectory **Middle**: Average trajectory length (# tokens) **Right**: Avg@4 scores on BrowseComp-en and BrowseComp-zh.

erQA, our model achieves the best performance among open-source sub-100B models. Specifically, with scores of 15.7% on BrowseComp-en, 32.0% on BrowseComp-zh, and 75.7% on FRAMES, our model surpasses these significantly larger models including WebSailor-72B and even larger model Kimi-K2. Similarly, on WebWalkerQA, we achieve 62.7%, outperforming WebShaper-72B by 10 absolute points. Our model also demonstrates competitive performance on additional benchmarks, achieving 50.0% on GAIA and 53.7% on XBench-DS, rivaling models with more parameters like WebThinker-32B and WebSailor-32B.

Strong Generalization from BrowseComp-Inspired Training Data Notably, while our QA pair synthesis methodology, particularly the evolution stage, is inspired by BrowseComp-en, and the exemplars provided in synthesis prompts are both derived from BrowseComp-en, we ensure diversity by collecting seed entities from Wikipedia across different domains. The resulting model also generalizes effectively across diverse benchmarks and domains beyond the BrowseComp-en/zh. As demonstrated previously, WEBEXPLORER-8B achieves strong performance across all information-seeking benchmarks, including GAIA, WebWalkerQA, FRAMES, and XBench-DS. This cross-benchmark success indicates substantial improvement in general information-seeking capabilities, demonstrating the generalization benefits of training on our challenging, complex queries. More remarkably, despite our training data not being STEM-focused, our model extends beyond information-seeking tasks. It achieves 17.3% on the HLE benchmark, which encompasses questions from diverse academic disciplines. This performance surpasses previous 32B models, like WebThinker-32B, further validating the robustness and transferability of our approach.

RL Training Dynamics and Performance Analysis To analyze the training dynamics, we recorded the average number of tool calls per trajectory and the average response length (number of tokens) during RL training, with results presented in Figure 4. Throughout the RL training process, the average number of tool calls increases steadily from approximately 11 to over 16, significantly exceeding the typical range of fewer than 10 tool calls observed in previous work (Gao et al., 2025). This substantial increase suggests that our model learns to execute more sophisticated multi-step reasoning strategies. Concurrently, the average trajectory length grows during the RL process, scaling to over 40K tokens, as the search tool responses constitute the majority source of trajectory length. Importantly, we observe that performance on both BrowseComp-en and BrowseComp-zh improves consistently throughout this process, with BrowseComp-en increasing from 7.9% to 15.7%, demonstrating a strong correlation between increased tool usage complexity and task performance. Similar to the emergence of long chain-of-thought phenomena observed in previous single-turn RL studies on mathematical or logical tasks (Guo et al., 2025; Zeng et al., 2025; Yu et al., 2025; Liu et al., 2025), these findings validate the effectiveness of our RL training approach in developing advanced web agents that can tackle complex information-seeking tasks through extended reasoning chains.

# 5 CONCLUSION

We present WEBEXPLORER, a simple framework for synthesizing high-quality information-seeking query-answer data for training web agents. Leveraging our data through SFT and RL, our WEB-EXPLORER-8B model achieves state-of-the-art performance across multiple information-seeking benchmarks, even outperforming larger models like WebSailor-72B on BrowseComp-en/zh. This success highlights a practical approach for building strong long-horizon web agents.

#### REPRODUCIBILITY STATEMENT

For our scaffold, we provide details about the used tools in Appendix D. For data synthesis, our WEBEXPLORER-QA framework is entirely prompting-based and designed for easy reproduction. We provide detailed prompting strategies and implementation guidelines in Appendix B. For the model, we use the widely adopted Qwen3-8B (Yang et al., 2025) and employ the GRPO (Shao et al., 2024) algorithm for reinforcement learning, ensuring our experimental setup can be readily replicated using standard tools and methodologies.

# REFERENCES

- Anthropic. Introducing claude 4. Technical report, 2025. URL https://www.anthropic.com/news/claude-4.
- Zijian Chen, Xueguang Ma, Shengyao Zhuang, Ping Nie, Kai Zou, Andrew Liu, Joshua Green, Kshama Patel, Ruoxi Meng, Mingyi Su, Sahel Sharifymoghaddam, Yanxi Li, Haoran Hong, Xinyu Shi, Xuye Liu, Nandan Thakur, Crystina Zhang, Luyu Gao, Wenhu Chen, and Jimmy Lin. Browsecomp-plus: A more fair and transparent evaluation benchmark of deep-research agent. arXiv preprint arXiv:2508.06600, 2025.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv* preprint arXiv:2507.06261, 2025.
- Jiaxuan Gao, Wei Fu, Minyang Xie, Shusheng Xu, Chuyi He, Zhiyu Mei, Banghua Zhu, and Yi Wu. Beyond ten turns: Unlocking long-horizon agentic search with large-scale asynchronous rl, 2025. URL https://arxiv.org/abs/2508.07976.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Jina.ai. Jina, 2025. URL https://jina.ai.
- Satyapriya Krishna, Kalpesh Krishna, Anhad Mohananey, Steven Schwarcz, Adam Stambler, Shyam Upadhyay, and Manaal Faruqui. Fact, fetch, and reason: A unified evaluation of retrieval-augmented generation. In *North American Chapter of the Association for Computational Linguistics*, 2024. URL https://api.semanticscholar.org/CorpusID:272753636.
- Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, et al. Websailor: Navigating super-human reasoning for web agent. *arXiv* preprint arXiv:2507.02592, 2025a.
- Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yutao Zhu, Yongkang Wu, Ji-Rong Wen, and Zhicheng Dou. Webthinker: Empowering large reasoning models with deep research capability. *arXiv preprint arXiv:2504.21776*, 2025b.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- Junteng Liu, Yuanxiang Fan, Zhuo Jiang, Han Ding, Yongyi Hu, Chi Zhang, Yiqi Shi, Shitong Weng, Aili Chen, Shiqi Chen, et al. Synlogic: Synthesizing verifiable reasoning data at scale for learning logical reasoning and beyond. *arXiv preprint arXiv:2505.19641*, 2025.
- Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. GAIA: a benchmark for general AI assistants. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=fibxvahvs3.
  - MiroMind Team. Miromind open deep research, 2025. URL https://miromind.ai/blog/miromind-open-deep-research. Blog post, August 2025.

- OpenAI. Introducing deep research. Technical report, 2025. URL https://openai.com/index/introducing-deep-research/.
- Vardaan Pahuja, Yadong Lu, Corby Rosset, Boyu Gou, Arindam Mitra, Spencer Whitehead, Yu Su, and Ahmed Hassan Awadallah. Explorer: Scaling exploration-driven web trajectory synthesis for multimodal web agents. In *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 6300–6323, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-256-5. URL https://aclanthology.org/2025.findings-acl.326/.
  - Perplexity AI. Introducing perplexity deep research. https://www.perplexity.ai/hub/blog/introducing-perplexity-deep-research, February 2025. Accessed: 2025-07-25.
    - Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, et al. Humanity's last exam. *arXiv preprint arXiv:2501.14249*, 2025.
    - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
    - Zhengwei Tao, Jialong Wu, Wenbiao Yin, Junkai Zhang, Baixuan Li, Haiyang Shen, Kuan Li, Liwen Zhang, Xinyu Wang, Yong Jiang, et al. Webshaper: Agentically data synthesizing via information-seeking formalization. *arXiv* preprint arXiv:2507.15061, 2025.
    - Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv* preprint arXiv:2507.20534, 2025.
    - Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet challenging benchmark for browsing agents. *arXiv preprint arXiv:2504.12516*, 2025.
    - Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Gang Fu, Yong Jiang, et al. Webdancer: Towards autonomous information seeking agency. *arXiv preprint arXiv:2505.22648*, 2025a.
    - Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Deyu Zhou, Pengjun Xie, and Fei Huang. Webwalker: Benchmarking llms in web traversal. *ArXiv*, abs/2501.07572, 2025b. URL https://api.semanticscholar.org/CorpusID: 275471576.
    - Xbench-Team. Xbench-deepsearch, 2025. URL https://xbench.org/agi/aisearch.
    - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
    - Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv* preprint arXiv:2503.18892, 2025.
  - Zhipu AI. GLM-4.5: Reasoning, Coding, and Agentic Abililties. https://z.ai/blog/glm-4.5, July 2025. Accessed: 2025-08-07.

Peilin Zhou, Bruce Leon, Xiang Ying, Can Zhang, Yifan Shao, Qichen Ye, Dading Chong, Zhiling Jin, Chenxuan Xie, Meng Cao, et al. Browsecomp-zh: Benchmarking web browsing ability of large language models in chinese. *arXiv preprint arXiv:2504.19314*, 2025.

# A TRAJECTORY EXAMPLE

 To illustrate how our agent framework operates in practice, we provide a detailed trajectory example from the WEBEXPLORER supervised fine-tuning data. This example demonstrates the agent's multistep reasoning process, thought, tool call, and tool response when handling complex queries.

**Question:** In a county adjacent to the nation's capital, a small press gaming company cofounded by someone whose previous pulp adventure game had won a silver-level award in a category focusing on game mechanics excellence. This person reached a specific age milestone during an online funding campaign spanning consecutive years that occurred exactly 6 years after their game's award. What was the exact final funding amount raised by this campaign, expressed in US dollars without the dollar sign or comma separators?

#### STEP 1

# Thought

#### <think>

This is a complex multi-part query that I need to break down systematically. Let me identify the key components:

1. A county adjacent to the nation's capital (so near Washington DC)

. . .

Let me start by searching for information about gaming companies near Washington DC that had crowdfunding campaigns, and awards for pulp adventure games.

</think>

I need to work through this step-by-step to identify the specific gaming company, award, and crowdfunding campaign. Let me start by searching for information about gaming companies near Washington DC and awards for pulp adventure games.

# Action

#### <tool\_call>

{"name": "search", "arguments": {"queries": ["small press gaming company Washington DC area", "pulp adventure game silver award game mechanics ENnie", "gaming company county adjacent Washington DC crowdfunding"]}}

</tool\_call>

# Observation

#### <tool\_response>

- search result for [small press gaming company Washington DC area] <title>op Washington, DC Gaming Companies (10) Built In </title>
- <url>https://builtin.com/companies/location/washington-dc/type/gaming-companies </url> <snippet>Top Washington, DC Gaming Companies (10) · Edge Systems Technologies · Some-
- thing Wicked Games · Polygon · Implicit Solution · BetMGM · Snaapiq · Microcoaching. </s-nippet>

</tool\_response>

648 649 650 STEP 2 651 652 Thought 653 654 <think> From the search results, I can see some relevant information: 655 1. I found mention of "Spirit of the Century" which won a Silver ENnie for Best Rules in 2007, 656 and it's described as "A really well put together game of pulp adventure" 657 658 </think> 659 Let me search more specifically for information about Spirit of the Century, which appears to be a pulp adventure game that won a silver ENnie, and its publisher. 660 661 662 663 <tool\_call> {"name": "search", "arguments": {"queries": ["Spirit of the Century silver ENnie 2007 Evil Hat Productions", "Evil Hat Productions Maryland Virginia crowdfunding Kickstarter", "Fred Hicks 666 Evil Hat Productions age milestone funding campaign"]}} 667 </tool\_call> 668 669 Observation 670 671 <tool\_response> 672 — search result for [Spirit of the Century silver ENnie 2007 Evil Hat Productions] — 673 <title>Spirit of the Century - Wikipedia </title> 674 <url>https://en.wikipedia.org/wiki/Spirit\_of\_the\_Century </url> 675 <snippet>Spirit of the Century is a pulp role-playing game published by Evil Hat Productions, 676 and based on Evil Hat's FATE system. ... 2007 Silver ENnie Award for Best ... </snippet> 677 </tool\_response> 679 680 STEP 3 681 Thought 683 684 685 Great! Now I have some key information: 686 1. Spirit of the Century is a pulp adventure game by Evil Hat Productions 687 688 Perfect! I found that Spirit of the Century by Evil Hat Productions won a Silver ENnie for Best 689 Rules in 2007, and Evil Hat is based in Silver Spring, Maryland (Montgomery County, adjacent 690 to DC). Now let me search for their crowdfunding campaigns around 2013 (6 years after 2007) 691 and information about Fred Hicks' age milestones. 692 693 Action 694 <tool\_call> 696 {"name": "search", "arguments": {"queries": ["Evil Hat Productions Kickstarter 2013 2014 697 Fred Hicks", "Fate Core Kickstarter 2013 funding amount final total", "Fred Hicks 40th birthday milestone age Fate Core campaign"]}} </tool\_call> 700

# Observation

#### <tool\_response>

- search result for [Evil Hat Productions Kickstarter 2013 2014 Fred Hicks] —
- <title>Race to Adventure! A Spirit of the Century<sup>TM</sup> Board Game Kickstarter </title>
- <url>https://www.kickstarter.com/projects/evilhat/race-to-adventure-a-spirit-of-the-centurytm-board/posts?ref=project\_share </url>
- <snippet>Evil Hat Productions is raising funds for Race to Adventure! A Spirit of the Century<sup>TM</sup> Board Game on Kickstarter! Race to Adventure!

</snippet>

</tool\_response>

#### STEP 4

#### Thought

#### <think>

I found some crucial information! The Fate Core Kickstarter by Evil Hat Productions raised exactly \$433,365 and ended on January 30, 2013. This is 6 years after the 2007 ENnie award for Spirit of the Century.

</think>

Excellent! I found that the Fate Core Kickstarter raised \$433,365 and ended on January 30, 2013 - which is exactly 6 years after the 2007 ENnie award. Let me get more details about this campaign to confirm it spans consecutive years and involves an age milestone.

#### Action

<tool\_call>

{"name": "browse", "arguments": {

"url": "https://www.kickstarter.com/projects/evilhat/fate-core",

"query": When did this campaign start and end? Did it span 2012 and 2013? Any mention of Fred Hicks birthday or age milestone during the campaign?}}

</tool\_call>

# Observation

#### <tool\_response>

The campaign started on \*\*December 3, 2012\*\*, and ended on \*\*January 29, 2013\*\*, spanning \*\*56 days\*\* and crossing over from \*\*2012 to 2013\*\*.

There is \*\*no mention of Fred Hicks' birthday or age milestone\*\* in the provided content.

</tool\_response>

# STEP 5

Thought <think> Perfect! This confirms that the Fate Core Kickstarter: </think> Perfect! The Fate Core campaign spanned consecutive years (December 3, 2012 to January 29, 2013) and occurred exactly 6 years after the 2007 ENnie. Now let me search for information about Fred Hicks' birthday and age milestones during this period. Action <tool\_call> {"name": "search", "arguments": {"queries": ["Fred Hicks birthday December January 40 years old", "Fred Hicks born 1972 1973 Evil Hat", "Fred Hicks 40th birthday December 2012 January 2013"]}} </tool\_call> Observation <tool\_response> — search result for [Fred Hicks birthday December January 40 years old] — <title>Happy birthday to the first general manager of Forsythe Associates ... </title> <url>https://www.facebook.com/groups/649530655242462/posts/1496564073872445 </url> <snippet>Happy birthday to the first general manager of Forsythe Associates [1962-1971], Fred Hicks [1920-1971]. </snippet> </tool\_response> STEP 6 Thought <think> Perfect! I found the crucial information. . . . </think> Therefore, the exact final funding amount raised by this campaign, expressed in US dollars with-out the dollar sign or comma separators, is: \*\*433365\*\* 

# B DETAILS OF WEBEXPLORER-QA SYNTHESIS FRAMEWORK

Compared with previous works (Li et al., 2025a; Tao et al., 2025) involving complex synthesis processes, our approach presents a simpler prompting-based framework. Here, we provide details about the synthesis methodology to facilitate reproducibility.

Our framework operates through two stages, each guided by carefully designed prompting strategies that leverage large language models for autonomous information gathering and query evolution.

#### B.1 MODEL-BASED EXPLORATION

We first collect a large number of entities from Wikipedia. In the first stage, we provide a seed entity as the search entry point along with three exemplar QA pairs from BrowseComp-en to demonstrate the desired question characteristics. The prompting strategy encourages the model to conduct search

and browsing activities starting from the seed entity, then synthesize a challenging query-answer pair using the collected knowledge.

The key instruction emphasizes creating challenging queries with subtle and obscured clues. We explicitly prompt the model to ensure that while the question should be challenging, the answer must remain unique and verifiable through the information space it constructs. This stage results in initial QA pairs that incorporate multi-website reasoning.

# Model-Based Exploration Prompt You need to create a challenging question for deep search based on real information. You should start by collecting information from the internet, then select a truth, and create a question where the truth needs to be discovered through search. You will start with a random "seed", then search and browse for whatever you want on the Internet, and create the question and truth from the information you gather. You should provide several subtle and blurred clues to make the question challenging, while ensuring the truth is unique. There are some examples: {examples} Let's start, with the seed of "{seed}". You need to provide the following information in the final <answer></answer> tag: <question> {{The challenging question you created based on real information.}} </question> <truth> {{The one and only exact truth to the question.}} </truth> IMPORTANT: You must include the <question> and <truth> tags in your final response for the system to parse your answer correctly. not provide any other response format.

### **B.2** Iterative Query Evolution

We provide the full trajectory from the first stage as input to the evolution process. The second stage systematically transforms the initial QA pair into a more challenging variant through strategic information reduction and obfuscation. The prompting strategy provides the model with the original question-answer pair and explicit instructions for three primary evolution mechanisms: (1) removing redundant or overly explicit descriptions that provide multiple pathways to the answer, (2) systematically obfuscating specific details such as dates, locations, and proper names with vaguer descriptors, and (3) searching for alternative terminologies to replace explicit entity mentions. The evolution process can iterate up to five times, with each cycle potentially increasing reasoning complexity.

This prompting-based approach eliminates the need for complex graph construction heuristics or predefined evolution rules, instead leveraging the model's natural language understanding and web exploration capabilities. The framework's simplicity enables easy and quick reproducibility for generating high-quality challenging QA pairs.

# **Iterative Query Evolution Prompt**

You need to make the following question more challenging while keeping the truth unique.

866 867

868

869

870

871

872

873

874

875

876 877

878

879

880

881

882

883

884

885

886 887

888

889

894 895 896

897

899

900

901

902903904905

906907908

909

910

915 916

917

Original question: {question} Original truth: {answer} You should make the question more challenging in the following wavs: 1. Remove some descriptions, especially when there are multiple descriptions that can lead to the truth answer in the question 2. Make one description more vague (such as date, location, name, etc.) in the question while keeping the truth answer unique 3. Search for new descriptions or alternative terms to replace specific entities in the question You can use search and browse tools in this process. Make sure the improved question is more challenging but the truth remains unique. You can iteratively make the question more and more challenging using these approaches up to 5 times. For each iteration, provide the evolved question within <question> and </question> tags. Provide the final improved question and truth pair in the final <answer></answer> tag: <question> {{The improved, more challenging and complex question.}} </question> <truth> {{The same exact truth.}} </truth> IMPORTANT: You must include the <question> and <truth> tags in your final response for the system to parse your answer correctly. not provide any other response format.

# C ILLUSTRATION OF MODEL-BASED EXPLORATION

Here, we provide an example to illustrate how the model-based exploration works in Figure 5. Starting from the seed "Brazil National Team", the model iteratively conducts search and browse actions to discover various interconnected facts. By focusing on one specific event—the "1976 FA Cup"—the model explores diverse information including details about players and referees, ultimately constructing a Query-Answer pair that requires deep reasoning across multiple discovered connections.

# D TOOL SCHEMAS

Our system utilizes two primary tools for web-based information retrieval: a search tool for discovering relevant content and a browse tool for detailed content analysis. The following sections detail the schemas for each tool.

#### D.1 SEARCH TOOL

The search tool enables web search functionality through Google search API, allowing the model to discover relevant URLs and content snippets based on input queries.

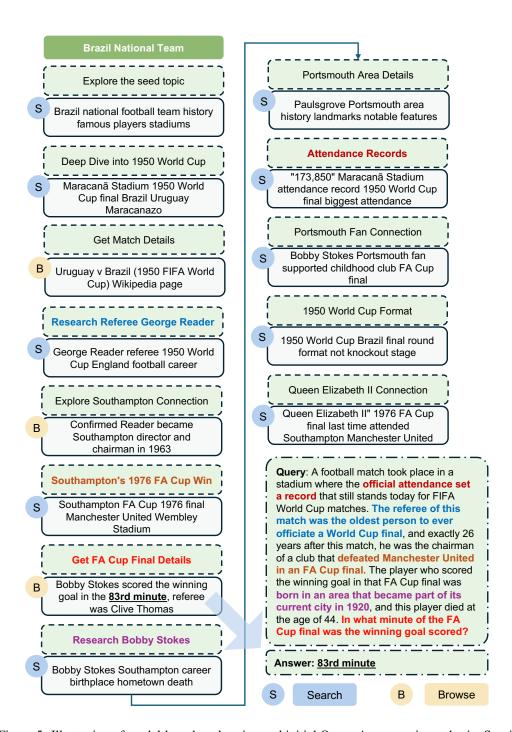


Figure 5: Illustration of model-based exploration and initial Query-Answer pair synthesis. Starting from the seed "Brazil National Team", the model iteratively explores using S (Search) and B (Browse) actions to discover interconnected facts, then synthesizes a challenging query-answer pair that requires deep reasoning across multiple discovered connections.

```
Search Tool Schema
tool_schema:
 type: function
  function:
   name: search
    description: Web search.
    parameters:
      type: object
      properties:
        queries:
          type: array
          description: The queries will be sent to Google. You will
                      get the brief search results with (title, url,
                      snippet)s for each query.
          items:
            type: string
      required:
        - queries
```

The search tool accepts an array of query strings and returns search results containing titles, URLs, and content snippets for each query. It is configured to return the top 10 results per query.

#### D.2 Browse Tool

972

973 974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989 990 991

992

993994995996

997 998

999

1000 1001

1023

1024

1025

The browse tool provides detailed content analysis capabilities by fetching and processing content from specific URLs using a combination of content extraction and language model analysis.

```
1002
          Browse Tool Schema
1003
1004
          tool schema:
1005
            type: function
            function:
              name: browse
1007
              description: Explore specific information in a url.
1008
              parameters:
1009
                type: object
1010
                properties:
                  url:
1011
                    type: string
1012
                     description: The url will be browsed, and the content
1013
                                 will be sent to a Large Language Model (LLM)
1014
                                  as the based information to answer a query.
1015
                  query:
1016
                    type: string
                     description: The query to this url content. You will
1017
                                 get an answer by another LLM.
1018
                required:
1019
                  - url
1020
                  - query
1021
1022
```

The browse tool takes a URL and a specific query as input. It uses Jina (Jina.ai, 2025) for content extraction and Gemini Flash (Comanici et al., 2025) as the generation engine to analyze the retrieved content and provide targeted answers to the input query.

# E THE USE OF LARGE LANGUAGE MODELS

We used large language models only for text polishing to improve grammar and readability. All intellectual contributions, including the experimental approach, analysis, and scientific insights, were developed solely by the authors.