LOST IN THE MAZE: OVERCOMING CONTEXT LIMITA-TIONS IN LONG-HORIZON INFORMATION SEEKING

Anonymous authors

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

020

021

024

025

026

027 028 029

031

032

033

034

037

038

040

041

043

044

046 047

048

051

052

Paper under double-blind review

ABSTRACT

Long-horizon information seeking tasks require iteratively searching the web over long trajectories and synthesizing information across many sources, and is a key capability to enable powerful applications like deep research systems. In this work, we show that popular information-seeking frameworks struggle to scale to long trajectories primarily due to context mismanagement—they accumulate long, noisy content, hit context window and tool budgets, or stop early. We introduce SLIM (Simple Lightweight Information Management), a simple framework that separates retrieval into distinct search and browse tools and periodically summarizes the trajectory, keeping context concise while enabling longer, more focused searches. On long-horizon tasks, SLIM achieves comparable accuracy at substantially lower cost and with far fewer tool calls than strong open-source baselines across multiple base models. Specifically, with o3 as the base model, SLIM achieves 55% on BrowseComp and 31% on HLE, outperforming all open-source frameworks by 7 and 3 absolute points, respectively, while incurring 5x fewer tool calls. Finally, we release an automated fine-grained trajectory analysis pipeline and error taxonomy for characterizing long-horizon information-seeking frameworks; SLIM exhibits less hallucination and fewer unfocused searches than prior systems. We hope our analysis framework and simple tool design inform future long-horizon agents.

1 Introduction

Long-horizon information-seeking tasks involve performing searches over long horizon and reasoning over many sources, and requires powerful systems that can explore diverse sources and leverage tools effectively. This core capability to reason over long-horizon serve as the foundation for exciting applications such as deep research (OpenAI, 2025; Google, 2025; xAI, 2025). Due to its immense potential in solving complex tasks, long-horizon systems been a key focus in the community, eliciting the development of many proprietary and open-source frameworks. For example, HuggingFace Open Deep Research (HF-ODR; Roucher et al., 2025) and GPT Researcher (GPT-R; Elovic, 2023) opts for complex multi-agent orchestration while SEARCH-O1 (Li et al., 2025b) uses a single-agent system. However, despite numerous implementations of these systems, they often fail to get the right answers in long-trajectory information-seeking settings and there are no systematic approaches to analyze their trajectories and identify the failure modes.

In this work, we first analyze existing frameworks by examining their trajectory outcomes on BrowseComp (Wei et al., 2025), a challenging long-horizon search benchmark. Our analysis shows that these frameworks still struggle with long-trajectory tasks, failing to achieve more than 50% accuracy—most of the failures are due to hitting either the context window limit, running out of tool budget, or stopping prematurely.

We attribute these failure modes to poor context management that can fill the context window often with noisy information that derails long search trajectories. The limited context restricts the number of turns in each trajectory, resulting in incomplete information gathering. To overcome these limitations, we design SLIM (Simple Lightweight Information Management), a framework with three simple yet powerful components—search, browse, and summarization—that effectively manage the context size of long-horizon systems. The simple tool design allows LLMs to interleave searching for diverse information and browsing promising pages without spending unnecessary tool calls on noisy search results. Furthermore, the summarization module acts as a general-purpose context manager that

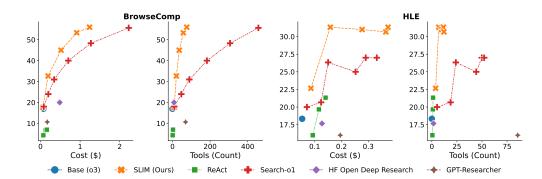


Figure 1: With o3 as the base model, SLIM achieves better performance than existing frameworks on both BrowseComp and HLE while using more than 5x fewer tool calls and lower overall costs, which accounts for LLM token usage and tool costs.

can reduce long trajectories into more condensed summaries. These design choices combine to allow the system to scale to longer trajectories while maintaining a concise context and reduced tool costs. Under a comparable cost budget, with o3 as the base model, SLIM significantly outperforms previous best open-source frameworks by 7 and 4 points on BrowseComp and HLE, respectively, while requiring only 25% of the tool calls (Figure 1).

Finally, we introduce an automated trajectory-level analysis pipeline that provides fine-grained insights to long-horizon frameworks. To characterize mistakes made by these systems, we develop an error taxonomy identifying common failure modes. Our analysis reveals that SLIM's advantage stems from its robustness to failure modes such as hallucinations and unfocused and generic searches. We hope our analysis pipeline, error taxonomy, and careful design choices in SLIM can serve as a foundation for understanding and improving long-horizon information-seeking systems.

2 Preliminaries: Long-Horizon Information Seeking

Previous information-seeking tasks, such as open-domain question answering, are relatively simple and straightforward as the queries are often factoid questions that are easy to answer with a single source of data (Joshi et al., 2017; Kwiatkowski et al., 2019; Petroni et al., 2021). As a result, these tasks can be mostly solved with static retrieval-augmented generation (RAG) systems that leverage at most a few retrieval steps (Lewis et al., 2020; Izacard et al., 2023; Shi et al., 2024), and do not showcase the challenges of realistic, long-horizon deep research settings.

In contrast, we study long-horizon tasks with complex queries that require extensive searches to gather the necessary information and reasoning over different sources to synthesize the answer. In this section, we formalize the task, describe the datasets for studying long-horizon information-seeking, and review some previous long-horizon systems.

2.1 TASK FORMULATION

We formalize long-horizon information-seeking tasks as follows: given a query q, a corpus of documents \mathcal{D} , the system needs to perform a sequence of tool calls to find relevant information from \mathcal{D} and output a final answer o, which is checked against the annotated groundtruth answer a. A critical component of the system is the design of its tools and how it interacts with the corpus; each tool is a function $\mathcal{T}_i(x) \to y$ that maps arbitrary system-generated inputs x to arbitrary outputs y.

Furthermore, agentic systems are often controlled by a tool budget T, which is the maximum number of tool calls they are allowed to use in any trajectory. The tool budget T also corresponds to the maximum number of turns in a trajectory, as each turn corresponds to one tool call or action. Thus, how to manage the input context to the underlying LLM across many tool uses and turns is another critical design choice in long-horizon systems.

¹Some architectures, such as the CodeAgent (Wang et al., 2024) used in HF-ODR, allow for parallel tool calls in one step, but we found that this does not occur in our long-horizon settings.

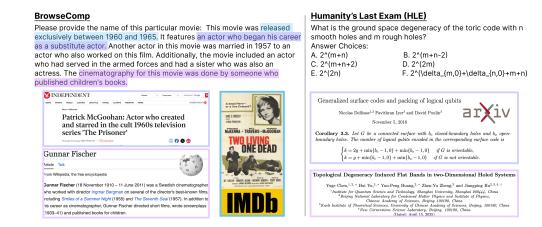


Figure 2: Example queries and their relevant documents for BrowseComp (Wei et al., 2025) and HLE (Phan et al., 2025).

In long-horizon information-seeking settings, \mathcal{D} is often set as the web due to the diversity and complexity of the queries, and each document $d_i = (u_i, t_i, c_i)$ is a tuple of its URL, title, and content. In practice, long-horizon systems typically use search engines $\mathcal{R}(q) \to \{(t_i, u_i)\}_1^n$ to obtain a list of n web pages with their titles and URLs most relevant to the search query q. Furthermore, a scraping operation $\mathcal{S}(u_i) \to c_i$ is necessary to obtain the full content of any URL as search engines only provide a list of URLs, but scraping is slow and noisy in practice.

In traditional QA settings, since the retrieval tool only needs to be called once due to the simplicity of the queries and the corpus (i.e., Wikipedia) is relatively small, the retrieval function returns the full list of documents and their contents $\mathcal{R}_{\text{wiki}}(q) \to \{(t_i, c_i)\}_1^n$. As a result, many existing long-horizon systems follow a similar design, where the retrieval tool is a single search engine call and scraping all returned documents. However, the complexity of long-horizon information-seeking tasks requires multiple rounds of tool calls to gather the necessary information. As we will demonstrate empirically later, this naive tool design leads to severe context limitations, where the system is overwhelmed by long, noisy content, motivating the design of more efficient tool interfaces for long-horizon systems.

2.2 Datasets

We select two datasets with naturally difficult queries that require long-trajectory searches. These datasets also have verifiable answers, which ensures the reliability of subsequent analyses. For evaluation, we sample a random subset of 300 instances from each dataset due to the expensive costs of running long-horizon systems. An example query from each dataset is shown in Figure 2.

BrowseComp (Wei et al., 2025) consists of challenging queries targeting hard-to-find information. BrowseComp tests one of the core capabilities of long-horizon systems—the ability to exhaustively search the web over long trajectories and collect the necessary information. These queries were rigorously validated by humans that even with access to the web, they cannot be solved within 10 minutes of searching. The answers are short and straightforward, resulting in reliable evaluation.

Humanity's Last Exam (HLE; Phan et al., 2025) tests across multiple domains and often require domain-specific knowledge and reasoning skills. These expert domains span across a wide range of topics, such as biology, mathematics, and physics. HLE tests the ability of long-horizon systems leverage the web to find helpful information that can aid reasoning-heavy problems. Similar to BrowseComp, the answers are easily verifiable due to their short lengths.

2.3 EXISTING APPROACHES

We briefly describe some popular approaches to information-seeking and deep research, ranging from simple single-LLM frameworks to complex multi-agent systems. We summarize the differences between these frameworks in Table 1; more details are in §A.1.

Table 1: Comparison of SLIM with existing frameworks. In contrast to previous works that bundle search and browsing search results into *one* retrieval tool, we separate it into two distinct tools.

Framework	Architecture	# Tools	Tools	Input to LLM Context	Summarization
REACT	Single-agent	1	Retrieval	All search results	-
SEARCH-01	Single-agent	1	Retrieval	All search results	Retrieved content
HF-ODR	Multi-agent	11	Search, Browse, Python,	Selected search results	Search agent result
GPT-R	Multi-agent	1	Retrieval	All search results	Retrieved content
SLIM (ours)	Single-agent	2	Search, Browse	Selected search results	Task trajectory

REACT (Yao et al., 2023) is a simple framework that allows an LLM agent to alternate between thinking and acting, allowing tool calling across many turns. Following the original work, our implementation gives the LLM access to a single retrieval tool—given a query, the tool returns a list of top 10 results along with their web contents. All search results are then concatenated and appended to the agent context for subsequent steps. When the LLM chooses not to use the search tool, the final output is used for evaluation. Our experiments vary the maximum number of turns in each trajectory.

SEARCH-01 (Li et al., 2025b) builds upon REACT with an additional "reason-in-document" step, where an LLM summarizes the search results and their contents before appending the results to the agent's input context. Although the summary step reduces context length for the main LLM compared to using full search result, this approach still uses a large amount of scrape operations in each search step, and summarization incurs a large amount of LLM token usage. We vary the maximum number of turns in each trajectory.

HuggingFace OpenDeepResearch (HF-ODR; Roucher et al., 2025) leverages a hierarchical structure consisting of a manager agent and a search agent. The manager agent calls the search agent to perform detailed searches, and the search agent iteratively interacts with a search engine and a browser, and returns a summary of its searches. The manager agent may use the summary to issue more queries or output the final answer. We use the default settings, which fixes the maximum number of turns for the manager and search agent to be 20.

GPT-Researcher (GPT-R; Elovic, 2023) is a complex multi-agent system where each agent has distinct roles: a research conductor that orchestrates the search process, a report generator that creates the report, a context manager that summarizes search results, and a source curator that selects relevant sources from scraped pages. The system uses a deep researcher agent that acts as a search tree node, spawning multiple children nodes with these same components. We use the default setting, which fixes the number of depths of the search tree to be 2 and the breath of search at each depth to be 4.

3 FAILURE MODES OF EXISTING APPROACHES

Despite recent progress, we still know little about how individual components in these systems perform, or fail. To study behavior on long-horizon tasks, we focus on BrowseComp, which naturally induces extended, multi-step search trajectories. For this task, the final outcome can reveal the overall performance of each framework as well as its relationship with the context window limitation and tool budget constraints. For this analysis, we let the framework run up to a fixed number of iterations and output an answer. We categorize the final outcome in Table 2.

Table 2: Categorization of different search outcomes and their descriptions

Outcome	Description
Correct	The system outputs the correct answer
Exceed context	The system exceeds LLM's context window, falling back to not using any tools
Exceed budget	The system exceeds the tool calling or iteration budget
Early stopping	The system outputs an incorrect answer before reaching the iteration budget
Misc. error	Due to uncontrollable factors (e.g., API content filters) the system outputs an error message

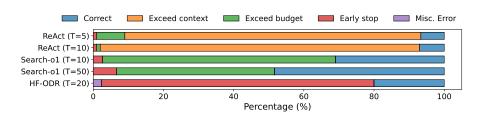


Figure 3: Each framework exhibit distinct outcome trends—REACT predominantly runs out of context window, while SEARCH-O1 is often limited by the tool budget (T). We exclude GPT-R due to its predefined workflow—the outcome can only be either correct or incorrect.

For this analysis, we consider different tool budgets for REACT and SEARCH-O1, and use the default 20 turns for HF-ODR, without loss of generality. We observe that context window limitations and tool budgets are the main bottlenecks for existing approaches in Figure 3, and each framework exhibits distinct patterns.

Specifically, REACT often hit the context window limit over a long trajectory due to the large amount of text returned by each search call. As a result, it cannot effectively scale to long trajectories and make full use of its tool budgets. SEARCH-O1 failure cases are almost entirely due to exceeding the tool budget, which suggests increasing the tool budget may potentially lead to better performance. However, such an increase is non-trivial without incurring a significant amount of cost—each retrieval step in SEARCH-O1 involves scraping all search results, even though only a fraction of these results are relevant, leading to a large amount of LLM token consumption during the summarization step.

Finally, we observe that HF-ODR often prematurely terminates due to the manager agent's inability to leverage its search agent across multiple steps. The root cause of these failure modes is bad context management—either exceeding context and tool budgets or stopping too early. In the next section we explore how to substantially improve deep research frameworks through better context management.

4 OUR FRAMEWORK: SLIM

A key takeaway from our analysis is that long-trajectory tasks require scaling up the number of turns and tool calls while keeping the context concise to avoid hitting the context window limit. Specifically, search results are often noisy and irrelevant to the answer, so filling up the context with content from all search results can lead to noisy context and unnecessary tool costs. Motivated by these observations, we introduce SLIM (Simple Lightweight Information Management) with two key principles: (1) using simple and flexible tools for LLMs to interact with, and (2) minimizing the amount of noisy information presented to the model and keeping the context concise during exploration. An overview of SLIM in comparison to existing frameworks are shown in Figure 4.

Concretely, SLIM adopts three simple yet powerful components—search, browse, and summarization—to effectively manage the context and scale the number of turns.

Search tool \mathcal{R} . As the main vehicle for exploring the web, SLIM uses a simple and fast interface for the search tool. Specifically, the search tool only returns the top k search results from a search engine, where each search result consists of a title, an URL, and a short snippet of its content. A crucial difference with previous frameworks is that previous work often bundles the search and browse functionality and return the full content for all search results, and relies on the main LLM to discern relevant context. In comparison, our search tool only returns a short snippet of each result, keeping the output concise and avoiding waste of context and tool use on irrelevant content.

Browse tool \mathcal{B} . Our browse tool is designed to complement the search tool by allowing the LLM to dig deeper into promising search results. Specifically, the browse tool $\mathcal{B}(u,q) \to \max_{c_i \in c} \text{sim}(c_i,q)$ returns the most relevant section of the content c from the URL u to the query q. Notably, this design enables the LLM to select the most relevant search result and choose a subset of the content that best matches the specific information it is looking for. As a result, our browse tool is significantly more

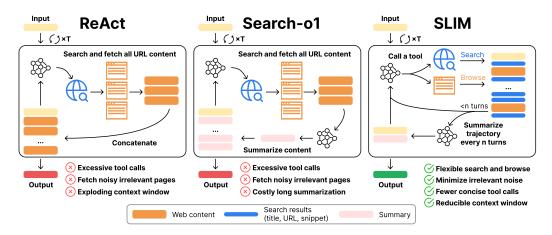


Figure 4: Compared to REACT and SEARCH-01, the cooperation between search, browse, and summarization modules allows SLIM to accumulate much shorter context length and less noisy information after exploring the same amount of searches.

efficient and cheaper than previous frameworks that exhaustively browse all search results in terms of both the scraping operations and the amount of new tokens introduced to the context.

Summarization module S. Despite the brevity of each tool response, agent context inevitably grows as it explores over a long horizon of searches. To maintain a concise context while retaining the effective exploration history, we introduce a summarization module that periodically compresses the LLM context. We find a simple heuristic sufficient: we summarize the entire conversation history after every n turns of tool calls and replace the trajectory with the summary. This crucially differs from previous works where summarization is solely applied to search results at each turn.

Finally, we combine these components into a single framework by allowing the underlying LLM to call either the search or the browse tools at every turn. Then, the summarization module compresses the entire conversation every n turns to reduce the amount of noise. Our implementation uses $Google^2$ as the search tool, $crawl4ai^3$ as the browse tool, and the same LLM as the agent model for summarization. More details and an example of the agent trajectory can be found in §A.4.

5 RESULTS

We use o3, o4-mini, and Claude-4-Sonnet as our base models. For each instance, we evaluate the system's performance as well as the number of tool calls and tokens used. The number of tool calls consider both the search API and browse/scraping operations. For the number of tokens, we take a weighted sum of the LLM input and output tokens across all turns. We exclude cached input tokens in total tokens count due to the fact that practical systems are implemented with caching mechanisms in long-trajectory tasks with many turns and shared context. For each dataset we record results averaged over all instances. More details on the experimental setup can be found in §A.5.

We present the main results with o3 as the base model in Table 3. Under the same cost, SLIM achieves significant improvements over SEARCH-O1, the best performing open-source framework, by 7 points on BrowseComp and 3 points on HLE. The difference is more pronounced when controlling for cost—SLIM can scale to 150 turns while using less total cost and reaching higher performance than SEARCH-O1 that can only scale to 50 turns. Furthermore, SLIM uses significantly fewer tool calls—less than 25% of the tool calls used by SEARCH-O1—suggesting that SLIM can leverage tools much more efficiently. The performance-cost comparisons of these systems are shown in Figure 1.

We also show results with different base models—o4-mini in Figure 5 and Claude-4-Sonnet in Figure 6. SLIM consistently achieves the highest performance across these models and all datasets compared to other frameworks, suggesting that our simple design generalizes well to models of

²https://serper.dev/

³https://github.com/unclecode/crawl4ai

Table 3: Main results with o3 as the base model. All results are macro-averaged across test instances. The number of tokens is shown in 10,000s. The cost is shown in US dollars. T denotes the maximum number of turns in each trajectory. For reference only, \dagger marks deep research systems that underwent task-specific training. Numbers are from the original reports (OpenAI, 2025; xAI, 2025; Qiao et al., 2025; Li et al., 2025c), and are not directly comparable due to different subsets of test instances used.

			Browse	Comp		HLE			
	T	Score (†)	Tokens (↓)	Tools (↓)	Cost (↓)	Score (†)	Tokens (↓)	Tools (↓)	Cost (↓)
o3 (no search)	-	17.0	3.8	0.0	0.08	18.3	2.7	0.0	0.05
	1	4.3	3.6	1.0	0.07	16.0	4.6	0.6	0.09
REACT	5	6.7	6.6	2.2	0.13	19.7	5.8	1.1	0.12
	10	7.0	8.0	2.8	0.16	21.3	7.0	1.2	0.14
	1	18.0	3.8	9.5	0.08	20.0	3.3	5.2	0.07
	5	24.0	8.0	46.9	0.20	20.7	5.4	18.7	0.12
SEARCH-01	10	31.0	13.7	89.8	0.35	26.3	6.6	23.9	0.15
SEARCH-OI	25	40.0	27.8	183.2	0.70	25.0	10.9	44.2	0.25
	50	48.3	51.5	306.2	1.27	27.0	12.6	49.8	0.29
	100	55.7	93.3	456.7	2.23	27.0	14.5	52.2	0.33
HF-ODR	20	20.0	24.1	8.4	0.49	17.7	6.4	1.7	0.13
GPT-R	-	10.7	5.8	69.5	0.17	16.0	6.4	85.6	0.20
	10	17.7	2.7	8.7	0.06	22.7	4.2	3.8	0.09
	25	32.7	9.0	20.7	0.19	31.3	7.7	6.9	0.16
SLIM	50	45.0	25.0	36.0	0.52	31.0	13.6	9.7	0.28
	100	53.3	44.1	57.4	0.91	31.3	18.4	11.6	0.37
	150	56.0	59.8	75.9	1.24	30.7	17.9	12.0	0.37
OpenAI DR [†]	-	51.5	-	-	-	26.6	-	-	-
Grok-4 [†]	-	44.9	-	_	-	38.6	-	-	-
WebS-30B [†]	-	37.3	-	-	-	28.8	-	-	-
WebT-32B [†]	-	15.8	-	-	-	-	-	-	-

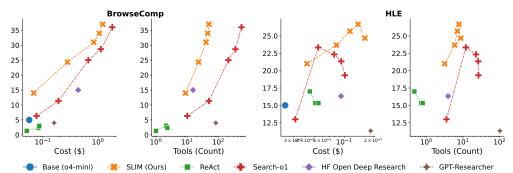


Figure 5: With o4-mini as the base model, SLIM consistently outperforms other baselines on BrowseComp while using fewer tool calls and lower overall costs. On HLE, SLIM can achieve overall higher performance and use fewer tool calls.

different sizes and training strategies. Furthermore, our effective context management also results in fewer tool calls and often lower overall costs compared to the baselines. SLIM also shows consistent trend across all three base models whereas certain framework only works well under certain models (e.g., HF-ODR only achieving competitive performance with Claude). Overall, this is strong evidence that SLIM serves as an effective framework for long-horizon tasks.

6 Fine-Grained Trajectory-level Analysis

6.1 Trajectory-Level Error Taxonomy

To understand how SLIM improves over other systems at a deeper level, we hope to go beyond analysis based solely on outcome, and focus on characterizing the mistakes that a system makes over the course of its long search *trajectories*. To this end, we first develop a shared taxonomy of common failure modes by manually examining individual trajectories from the compared systems on BrowseComp.

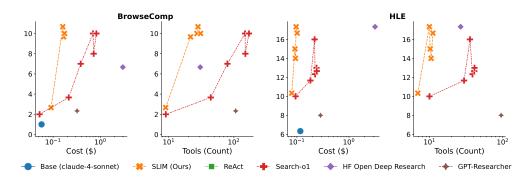


Figure 6: With Claude-4-Sonnet as the base model, SLIM consistently outperforms other baselines on BrowseComp while using fewer tool calls and lower overall costs. On HLE, SLIM can achieve overall higher performance and use fewer tool calls.

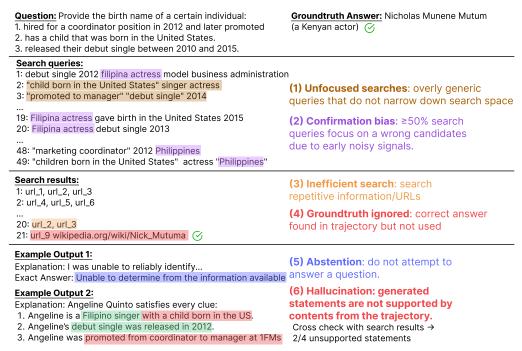


Figure 7: Examples of each trajectory-level failure mode on a BrowseComp sample.

Then based on the taxonomy, we develop an automated error analysis pipeline that annotates each trajectory with the failure modes using a mix of rule-based heuristics and LLM-as-a-judge approaches.

We present examples of each failure mode in the taxonomy in Figure 7. We also include detailed definitions of these failure modes and how we construct the error analysis pipeline in §A.3.

6.2 Analysis of Trajectory-Level Failure Modes

For fair comparison, we analyze all frameworks under a similar cost budget. For each framework we choose the setting with closest cost to SLIM with 150 iterations, according to Table 3. The distribution of trajectory-level errors are shown in Table 4. We first observe that SLIM's advantage in performance could be attributed to the notably reduced number of unfocused searches and hallucination rate. Since SLIM can issue more search queries due to its relatively low cost compared to other frameworks that consume longer context for search and browsing, it is able to conduct wider and more diverse searches, with only 37.0% trajectories exhibiting unfocused searches whereas other frameworks observe this in more than 40% of their trajectories.

Table 4: The percentage of trajectory that observe each failure mode.

Framework	Turn Budget	Correct	Confirm Bias	Unfocused Search	Inefficient Search	Abstention	Answer Ignored	Hallucinate
REACT	10	7.0	32.7	54.0	3.9	1.0	0.7	40.0
SEARCH-01	50	48.3	21.0	42.0	7.2	4.3	23.0	45.2
HF-ODR	20	20.0	15.7	44.3	43.9	32.0	0.0	95.9
GPT-R	-	10.7	25.7	73.3	-	12.3	-	-
SLIM	150	56.0	19.7	37.0	7.6	27.3	27.3	19.1

Furthermore, SLIM carefully manages the amount of information added to the context and minimizes the amount of noise, leading to an overall lower hallucination rate of only 14.5% compared to the next lowest of 40% from REACT. In contrast, frameworks that consume significantly more cost at issuing search queries—such as HF-ODR—are more prone to hallucination since they often resort to their parametric knowledge to answer the question.

Notably, despite the improvements on these two axes and an overall higher correct ratio, SLIM still suffers from high confirmation bias and abstention, and is more prone to ignoring the groundtruth answers. We leave these improvements to future work, and hope that our trajectory-level analysis can be a useful tool for improving long-horizon systems in more interpretable and concrete ways.

7 RELATED WORK

 Deep research. Recently, the community has taken great interests in deep research systems due to their potential to solve complex tasks—there have been several efforts across both industry (OpenAI, 2025; Google, 2025; xAI, 2025; Nguyen et al., 2025) and open-source communities (Wu et al., 2025a; Du et al., 2025; Sun et al., 2025, *inter alia*). They are often evaluated through long-horizon search trajectories tasks that also require complex reasoning (Wei et al., 2025; Phan et al., 2025). Other benchmarks evaluate the long-form generation capabilities of systems (Du et al., 2025).

Furthermore, between the opaque proprietary systems and increasingly complex open-source systems, there is little understanding on the underlying behavior of long-horizon systems and how they fail in practice. In this work, we aim to fill this gap by introducing a error taxonomy for long-horizon systems and an automatic error analysis pipeline. Finally, in contrast to existing open-source approaches that are growing increasingly more complex, we show that a simple approach with carefully designed tools can achieve better performance with fewer tool calls.

Reinforcement learning for long-horizon systems. There have been considerable efforts in improving search agents through reinforcement learning (Li et al., 2025c; Zheng et al., 2025; Chen et al., 2025; Li et al., 2025a; Wu et al., 2025b, *inter alia*). A popular approach is to synthetically generate question-answer pairs that require long-horizon search trajectories (Xia et al., 2025; Tao et al., 2025). Other works focuses on comparing different training objectives (Jin et al., 2025b;a). However, critical analysis of the error modes and comparison of different frameworks are still lacking.

8 CONCLUSION

In this work, we propose SLIM, a simple yet effective long-horizon information-seeking framework that address context limitations prevalent in existing systems. We show that SLIM consistently achieves the highest performance across different base models and datasets compared to other frameworks while using fewer tool calls and lower overall costs, suggesting that our simple enables better long-horizon information-seeking.

We then develop an automated error analysis pipeline to characterize the failure modes of long-horizon systems. Our analysis shows that SLIM is more resistant to failure modes such as unfocused search and hallucination. We hope our framework and analysis pipeline can serve as a useful tool for the community to understand and improve long-horizon systems.

ETHICS STATEMENT

This work studies the behavior of long-horizon information-seeking systems, and how to improve them through better design choices. Although there are no direct ethical concerns, we acknowledge that the web and LLMs are complex systems that can be used for harmful purposes.

REPRODUCIBILITY STATEMENT

We submit our code and data in the Supplemental Material, which includes everything needed to reproduce our results. We plan to release the code and data publicly after the review process.

REFERENCES

- Bernd Bohnet, Vinh Q Tran, Pat Verga, Roee Aharoni, Daniel Andor, Livio Baldini Soares, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, Kai Hui, et al. Attributed question answering: Evaluation and modeling for attributed large language models. *arXiv preprint arXiv:2212.08037*, 2022. URL https://arxiv.org/pdf/2212.08037.pdf.
- Mingyang Chen, Linzhuang Sun, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Haofen Wang, Jeff Z. Pan, Wen Zhang, Huajun Chen, Fan Yang, Zenan Zhou, and Weipeng Chen. Research: Learning to reason with search for llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2503.19470.
- Mingxuan Du, Benfeng Xu, Chiwei Zhu, Xiaorui Wang, and Zhendong Mao. Deepresearch bench: A comprehensive benchmark for deep research agents, 2025. URL https://arxiv.org/abs/2506.11763.
- Assaf Elovic. gpt-researcher, July 2023. URL https://github.com/assafelovic/gpt-researcher.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. Enabling large language models to generate text with citations. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2023.
- Google. Gemini deep research your personal research assistant, September 2025. URL https://gemini.google/overview/deep-research/.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: few-shot learning with retrieval augmented language models. *J. Mach. Learn. Res.*, 24(1), January 2023. ISSN 1532-4435.
- Bowen Jin, Jinsung Yoon, Priyanka Kargupta, Sercan O. Arik, and Jiawei Han. An empirical study on reinforcement learning for reasoning-search interleaved llm agents, 2025a. URL https://arxiv.org/abs/2505.15117.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan O Arik, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training LLMs to reason and leverage search engines with reinforcement learning. In *Second Conference on Language Modeling*, 2025b. URL https://openreview.net/forum?id=Rwhi9lideu.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1601–1611, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1147. URL https://aclanthology.org/P17-1147.
- Ryo Kamoi, Tanya Goyal, Juan Diego Rodriguez, and Greg Durrett. WiCE: Real-World Entailment for Claims in Wikipedia. *arXiv preprint arXiv:2303.01432*, 2023. URL https://arxiv.org/abs/2303.01432.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019. doi: 10.1162/tacl_a_00276. URL https://aclanthology.org/Q19-1026.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.

Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, Weizhou Shen, Junkai Zhang, Dingchu Zhang, Xixi Wu, Yong Jiang, Ming Yan, Pengjun Xie, Fei Huang, and Jingren Zhou. Websailor: Navigating super-human reasoning for web agent, 2025a. URL https://arxiv.org/abs/2507.02592.

Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. Search-ol: Agentic search-enhanced large reasoning models, 2025b. URL https://arxiv.org/abs/2501.05366.

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. *CoRR*, abs/2504.21776, 2025c. doi: 10.48550/ARXIV.2504.21776. URL https://doi.org/10.48550/arXiv.2504.21776.

Xuan-Phi Nguyen, Shrey Pandit, Revanth Gangi Reddy, Austin Xu, Silvio Savarese, Caiming Xiong, and Shafiq Joty. Sfr-deepresearch: Towards effective reinforcement learning for autonomously reasoning single agents, 2025. URL https://arxiv.org/abs/2509.06283.

OpenAI. Introducing deep research, February 2025. URL https://openai.com/index/ introducing-deep-research/.

Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. KILT: a benchmark for knowledge intensive language tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2523–2544, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.200. URL https://aclanthology.org/2021.naacl-main.200.

Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, Michael Choi, Anish Agrawal, Arnav Chopra, Adam Khoja, Ryan Kim, Richard Ren, Jason Hausenloy, Oliver Zhang, Mantas Mazeika, Dmitry Dodonov, Tung Nguyen, Jaeho Lee, Daron Anderson, Mikhail Doroshenko, Alun Cennyth Stokes, Mobeen Mahmood, Oleksandr Pokutnyi, Oleg Iskra, Jessica P. Wang, John-Clark Levin, Mstyslav Kazakov, Fiona Feng, Steven Y. Feng, Haoran Zhao, Michael Yu, Varun Gangal, Chelsea Zou, Zihan Wang, Serguei Popov, Robert Gerbicz, Geoff Galgon, Johannes Schmitt, Will Yeadon, Yongki Lee, Scott Sauers, Alvaro Sanchez, Fabian Giska, Marc Roth, Søren Riis, Saiteja Utpala, Noah Burns, Gashaw M. Goshu, Mohinder Maheshbhai Naiya, Chidozie Agu, Zachary Giboney, Antrell Cheatom, Francesco Fournier-Facio, Sarah-Jane Crowson, Lennart Finke, Zerui Cheng, Jennifer Zampese, Ryan G. Hoerr, Mark Nandor, Hyunwoo Park, Tim Gehrunger, Jiaqi Cai, Ben McCarty, Alexis C Garretson, Edwin Taylor, Damien Sileo, Qiuyu Ren, Usman Qazi, Lianghui Li, Jungbae Nam, John B. Wydallis, Pavel Arkhipov, Jack Wei Lun Shi, Aras Bacho, Chris G. Willcocks, Hangrui Cao, Sumeet Motwani, Emily de Oliveira Santos, Johannes Veith, Edward Vendrow, Doru Cojoc, Kengo Zenitani, Joshua Robinson, Longke Tang, Yuqi Li, Joshua Vendrow, Natanael Wildner Fraga, Vladyslav Kuchkin, Andrey Pupasov Maksimov, Pierre Marion, Denis Efremov, Jayson Lynch, Kaiqu Liang, Aleksandar Mikov, Andrew Gritsevskiy, Julien Guillod, Gözdenur Demir, Dakotah Martinez, Ben Pageler, Kevin Zhou, Saeed Soori, Ori Press, Henry Tang, Paolo Rissone, Sean R. Green, Lina Brüssel, Moon Twayana, Aymeric Dieuleveut, Joseph Marvin

595

596

597

598

600

601

602

603

604

605

606

607

608

610

611

612

613

614

615

616

617

618

619

620

621

622

623

625

626

627

629

630

631

632

633

634

635

636

637

638

639

640

641

642

644

645

646

Imperial, Ameya Prabhu, Jinzhou Yang, Nick Crispino, Arun Rao, Dimitri Zvonkine, Gabriel Loiseau, Mikhail Kalinin, Marco Lukas, Ciprian Manolescu, Nate Stambaugh, Subrata Mishra, Tad Hogg, Carlo Bosio, Brian P Coppola, Julian Salazar, Jaehyeok Jin, Rafael Sayous, Stefan Ivanov, Philippe Schwaller, Shaipranesh Senthilkuma, Andres M Bran, Andres Algaba, Kelsey Van den Houte, Lynn Van Der Sypt, Brecht Verbeken, David Noever, Alexei Kopylov, Benjamin Myklebust, Bikun Li, Lisa Schut, Evgenii Zheltonozhskii, Qiaochu Yuan, Derek Lim, Richard Stanley, Tong Yang, John Maar, Julian Wykowski, Martí Oller, Anmol Sahu, Cesare Giulio Ardito, Yuzheng Hu, Ariel Ghislain Kemogne Kamdoum, Alvin Jin, Tobias Garcia Vilchis, Yuexuan Zu, Martin Lackner, James Koppel, Gongbo Sun, Daniil S. Antonenko, Steffi Chern, Bingchen Zhao, Pierrot Arsene, Joseph M Cavanagh, Daofeng Li, Jiawei Shen, Donato Crisostomi, Wenjin Zhang, Ali Dehghan, Sergey Ivanov, David Perrella, Nurdin Kaparov, Allen Zang, Ilia Sucholutsky, Arina Kharlamova, Daniil Orel, Vladislav Poritski, Shalev Ben-David, Zachary Berger, Parker Whitfill, Michael Foster, Daniel Munro, Linh Ho, Shankar Sivarajan, Dan Bar Hava, Aleksey Kuchkin, David Holmes, Alexandra Rodriguez-Romero, Frank Sommerhage, Anji Zhang, Richard Moat, Keith Schneider, Zakayo Kazibwe, Don Clarke, Dae Hyun Kim, Felipe Meneguitti Dias, Sara Fish, Veit Elser, Tobias Kreiman, Victor Efren Guadarrama Vilchis, Immo Klose, Ujjwala Anantheswaran, Adam Zweiger, Kaivalya Rawal, Jeffery Li, Jeremy Nguyen, Nicolas Daans, Haline Heidinger, Maksim Radionov, Václav Rozhoň, Vincent Ginis, Christian Stump, Niv Cohen, Rafał Poświata, Josef Tkadlec, Alan Goldfarb, Chenguang Wang, Piotr Padlewski, Stanislaw Barzowski, Kyle Montgomery, Ryan Stendall, Jamie Tucker-Foltz, Jack Stade, T. Ryan Rogers, Tom Goertzen, Declan Grabb, Abhishek Shukla, Alan Givré, John Arnold Ambay, Archan Sen, Muhammad Fayez Aziz, Mark H Inlow, Hao He, Ling Zhang, Younesse Kaddar, Ivar Ängquist, Yanxu Chen, Harrison K Wang, Kalyan Ramakrishnan, Elliott Thornley, Antonio Terpin, Hailey Schoelkopf, Eric Zheng, Avishy Carmi, Ethan D. L. Brown, Kelin Zhu, Max Bartolo, Richard Wheeler, Martin Stehberger, Peter Bradshaw, JP Heimonen, Kaustubh Sridhar, Ido Akov, Jennifer Sandlin, Yury Makarychev, Joanna Tam, Hieu Hoang, David M. Cunningham, Vladimir Goryachev, Demosthenes Patramanis, Michael Krause, Andrew Redenti, David Aldous, Jesyin Lai, Shannon Coleman, Jiangnan Xu, Sangwon Lee, Ilias Magoulas, Sandy Zhao, Ning Tang, Michael K. Cohen, Orr Paradise, Jan Hendrik Kirchner, Maksym Ovchynnikov, Jason O. Matos, Adithya Shenoy, Michael Wang, Yuzhou Nie, Anna Sztyber-Betley, Paolo Faraboschi, Robin Riblet, Jonathan Crozier, Shiv Halasyamani, Shreyas Verma, Prashant Joshi, Eli Meril, Ziqiao Ma, Jérémy Andréoletti, Raghav Singhal, Jacob Platnick, Volodymyr Nevirkovets, Luke Basler, Alexander Ivanov, Seri Khoury, Nils Gustafsson, Marco Piccardo, Hamid Mostaghimi, Qijia Chen, Virendra Singh, Tran Quoc Khánh, Paul Rosu, Hannah Szlyk, Zachary Brown, Himanshu Narayan, Aline Menezes, Jonathan Roberts, William Alley, Kunyang Sun, Arkil Patel, Max Lamparth, Anka Reuel, Linwei Xin, Hanmeng Xu, Jacob Loader, Freddie Martin, Zixuan Wang, Andrea Achilleos, Thomas Preu, Tomek Korbak, Ida Bosio, Fereshteh Kazemi, Zive Chen, Biró Bálint, Eve J. Y. Lo, Jiaqi Wang, Maria Inês S. Nunes, Jeremiah Milbauer, M Saiful Bari, Zihao Wang, Behzad Ansarinejad, Yewen Sun, Stephane Durand, Hossam Elgnainy, Guillaume Douville, Daniel Tordera, George Balabanian, Hew Wolff, Lynna Kvistad, Hsiaoyun Milliron, Ahmad Sakor, Murat Eron, Andrew Favre D. O., Shailesh Shah, Xiaoxiang Zhou, Firuz Kamalov, Sherwin Abdoli, Tim Santens, Shaul Barkan, Allison Tee, Robin Zhang, Alessandro Tomasiello, G. Bruno De Luca, Shi-Zhuo Looi, Vinh-Kha Le, Noam Kolt, Jiayi Pan, Emma Rodman, Jacob Drori, Carl J Fossum, Niklas Muennighoff, Milind Jagota, Ronak Pradeep, Honglu Fan, Jonathan Eicher, Michael Chen, Kushal Thaman, William Merrill, Moritz Firsching, Carter Harris, Stefan Ciobâcă, Jason Gross, Rohan Pandey, Ilya Gusev, Adam Jones, Shashank Agnihotri, Pavel Zhelnov, Mohammadreza Mofayezi, Alexander Piperski, David K. Zhang, Kostiantyn Dobarskyi, Roman Leventov, Ignat Soroko, Joshua Duersch, Vage Taamazyan, Andrew Ho, Wenjie Ma, William Held, Ruicheng Xian, Armel Randy Zebaze, Mohanad Mohamed, Julian Noah Leser, Michelle X Yuan, Laila Yacar, Johannes Lengler, Katarzyna Olszewska, Claudio Di Fratta, Edson Oliveira, Joseph W. Jackson, Andy Zou, Muthu Chidambaram, Timothy Manik, Hector Haffenden, Dashiell Stander, Ali Dasouqi, Alexander Shen, Bita Golshani, David Stap, Egor Kretov, Mikalai Uzhou, Alina Borisovna Zhidkovskaya, Nick Winter, Miguel Orbegozo Rodriguez, Robert Lauff, Dustin Wehr, Colin Tang, Zaki Hossain, Shaun Phillips, Fortuna Samuele, Fredrik Ekström, Angela Hammon, Oam Patel, Faraz Farhidi, George Medley, Forough Mohammadzadeh, Madellene Peñaflor, Haile Kassahun, Alena Friedrich, Rayner Hernandez Perez, Daniel Pyda, Taom Sakal, Omkar Dhamane, Ali Khajegili Mirabadi, Eric Hallman, Kenchi Okutsu, Mike Battaglia, Mohammad Maghsoudimehrabani, Alon Amit, Dave Hulbert, Roberto Pereira, Simon Weber, Handoko, Anton Peristyy, Stephen Malina, Mustafa Mehkary, Rami Aly, Frank Reidegeld, Anna-Katharina Dick, Cary Friday, Mukhwinder Singh, Hassan Shapourian, Wanyoung Kim, Mar-

650

651

652

653

654

655

656

657

658

659

660

661

662

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

687

688

689

690

691

692

693

694

696

699

700

iana Costa, Hubeyb Gurdogan, Harsh Kumar, Chiara Ceconello, Chao Zhuang, Haon Park, Micah Carroll, Andrew R. Tawfeek, Stefan Steinerberger, Daattavya Aggarwal, Michael Kirchhof, Linjie Dai, Evan Kim, Johan Ferret, Jainam Shah, Yuzhou Wang, Minghao Yan, Krzysztof Burdzy, Lixin Zhang, Antonio Franca, Diana T. Pham, Kang Yong Loh, Joshua Robinson, Abram Jackson, Paolo Giordano, Philipp Petersen, Adrian Cosma, Jesus Colino, Colin White, Jacob Votava, Vladimir Vinnikov, Ethan Delaney, Petr Spelda, Vit Stritecky, Syed M. Shahid, Jean-Christophe Mourrat, Lavr Vetoshkin, Koen Sponselee, Renas Bacho, Zheng-Xin Yong, Florencia de la Rosa, Nathan Cho, Xiuyu Li, Guillaume Malod, Orion Weller, Guglielmo Albani, Leon Lang, Julien Laurendeau, Dmitry Kazakov, Fatimah Adesanya, Julien Portier, Lawrence Hollom, Victor Souza, Yuchen Anna Zhou, Julien Degorre, Yiğit Yalın, Gbenga Daniel Obikoya, Rai, Filippo Bigi, M. C. Boscá, Oleg Shumar, Kaniuar Bacho, Gabriel Recchia, Mara Popescu, Nikita Shulga, Ngefor Mildred Tanwie, Thomas C. H. Lux, Ben Rank, Colin Ni, Matthew Brooks, Alesia Yakimchyk, Huanxu, Liu, Stefano Cavalleri, Olle Häggström, Emil Verkama, Joshua Newbould, Hans Gundlach, Leonor Brito-Santana, Brian Amaro, Vivek Vajipey, Rynaa Grover, Ting Wang, Yosi Kratish, Wen-Ding Li, Sivakanth Gopi, Andrea Caciolai, Christian Schroeder de Witt, Pablo Hernández-Cámara, Emanuele Rodolà, Jules Robins, Dominic Williamson, Vincent Cheng, Brad Raynor, Hao Qi, Ben Segev, Jingxuan Fan, Sarah Martinson, Erik Y. Wang, Kaylie Hausknecht, Michael P. Brenner, Mao Mao, Christoph Demian, Peyman Kassani, Xinyu Zhang, David Avagian, Eshawn Jessica Scipio, Alon Ragoler, Justin Tan, Blake Sims, Rebeka Plecnik, Aaron Kirtland, Omer Faruk Bodur, D. P. Shinde, Yan Carlos Leyva Labrador, Zahra Adoul, Mohamed Zekry, Ali Karakoc, Tania C. B. Santos, Samir Shamseldeen, Loukmane Karim, Anna Liakhovitskaia, Nate Resman, Nicholas Farina, Juan Carlos Gonzalez, Gabe Maayan, Earth Anderson, Rodrigo De Oliveira Pena, Elizabeth Kelley, Hodjat Mariji, Rasoul Pouriamanesh, Wentao Wu, Ross Finocchio, Ismail Alarab, Joshua Cole, Danyelle Ferreira, Bryan Johnson, Mohammad Safdari, Liangti Dai, Siriphan Arthornthurasuk, Isaac C. McAlister, Alejandro José Moyano, Alexey Pronin, Jing Fan, Angel Ramirez-Trinidad, Yana Malysheva, Daphiny Pottmaier, Omid Taheri, Stanley Stepanic, Samuel Perry, Luke Askew, Raúl Adrián Huerta Rodríguez, Ali M. R. Minissi, Ricardo Lorena, Krishnamurthy Iyer, Arshad Anil Fasiludeen, Ronald Clark, Josh Ducey, Matheus Piza, Maja Somrak, Eric Vergo, Juehang Qin, Benjámin Borbás, Eric Chu, Jack Lindsey, Antoine Jallon, I. M. J. McInnis, Evan Chen, Avi Semler, Luk Gloor, Tej Shah, Marc Carauleanu, Pascal Lauer, Tran Duc Huy, Hossein Shahrtash, Emilien Duc, Lukas Lewark, Assaf Brown, Samuel Albanie, Brian Weber, Warren S. Vaz, Pierre Clavier, Yiyang Fan, Gabriel Poesia Reis e Silva, Long, Lian, Marcus Abramovitch, Xi Jiang, Sandra Mendoza, Murat Islam, Juan Gonzalez, Vasilios Mavroudis, Justin Xu, Pawan Kumar, Laxman Prasad Goswami, Daniel Bugas, Nasser Heydari, Ferenc Jeanplong, Thorben Jansen, Antonella Pinto, Archimedes Apronti, Abdallah Galal, Ng Ze-An, Ankit Singh, Tong Jiang, Joan of Arc Xavier, Kanu Priya Agarwal, Mohammed Berkani, Gang Zhang, Zhehang Du, Benedito Alves de Oliveira Junior, Dmitry Malishev, Nicolas Remy, Taylor D. Hartman, Tim Tarver, Stephen Mensah, Gautier Abou Loume, Wiktor Morak, Farzad Habibi, Sarah Hoback, Will Cai, Javier Gimenez, Roselynn Grace Montecillo, Jakub Łucki, Russell Campbell, Asankhaya Sharma, Khalida Meer, Shreen Gul, Daniel Espinosa Gonzalez, Xavier Alapont, Alex Hoover, Gunjan Chhablani, Freddie Vargus, Arunim Agarwal, Yibo Jiang, Deepakkumar Patil, David Outevsky, Kevin Joseph Scaria, Rajat Maheshwari, Abdelkader Dendane, Priti Shukla, Ashley Cartwright, Sergei Bogdanov, Niels Mündler, Sören Möller, Luca Arnaboldi, Kunvar Thaman, Muhammad Rehan Siddiqi, Prajvi Saxena, Himanshu Gupta, Tony Fruhauff, Glen Sherman, Mátyás Vincze, Siranut Usawasutsakorn, Dylan Ler, Anil Radhakrishnan, Innocent Enyekwe, Sk Md Salauddin, Jiang Muzhen, Aleksandr Maksapetyan, Vivien Rossbach, Chris Harjadi, Mohsen Bahaloohoreh, Claire Sparrow, Jasdeep Sidhu, Sam Ali, Song Bian, John Lai, Eric Singer, Justine Leon Uro, Greg Bateman, Mohamed Sayed, Ahmed Menshawy, Darling Duclosel, Dario Bezzi, Yashaswini Jain, Ashley Aaron, Murat Tiryakioglu, Sheeshram Siddh, Keith Krenek, Imad Ali Shah, Jun Jin, Scott Creighton, Denis Peskoff, Zienab EL-Wasif, Ragavendran P V, Michael Richmond, Joseph McGowan, Tejal Patwardhan, Hao-Yu Sun, Ting Sun, Nikola Zubić, Samuele Sala, Stephen Ebert, Jean Kaddour, Manuel Schottdorf, Dianzhuo Wang, Gerol Petruzella, Alex Meiburg, Tilen Medved, Ali ElSheikh, S Ashwin Hebbar, Lorenzo Vaquero, Xianjun Yang, Jason Poulos, Vilém Zouhar, Sergey Bogdanik, Mingfang Zhang, Jorge Sanz-Ros, David Anugraha, Yinwei Dai, Anh N. Nhu, Xue Wang, Ali Anil Demircali, Zhibai Jia, Yuyin Zhou, Juncheng Wu, Mike He, Nitin Chandok, Aarush Sinha, Gaoxiang Luo, Long Le, Mickaël Noyé, Michał Perełkiewicz, Ioannis Pantidis, Tianbo Qi, Soham Sachin Purohit, Letitia Parcalabescu, Thai-Hoa Nguyen, Genta Indra Winata, Edoardo M. Ponti, Hanchen Li, Kaustubh Dhole, Jongee Park, Dario Abbondanza, Yuanli Wang, Anupam Nayak, Diogo M. Caetano, Antonio A. W. L. Wong, Maria del Rio-Chanona, Dániel

703

704

705

706

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

Kondor, Pieter Francois, Ed Chalstrey, Jakob Zsambok, Dan Hoyer, Jenny Reddish, Jakob Hauser, Francisco-Javier Rodrigo-Ginés, Suchandra Datta, Maxwell Shepherd, Thom Kamphuis, Qizheng Zhang, Hyunjun Kim, Ruiji Sun, Jianzhu Yao, Franck Dernoncourt, Satyapriya Krishna, Sina Rismanchian, Bonan Pu, Francesco Pinto, Yingheng Wang, Kumar Shridhar, Kalon J. Overholt, Glib Briia, Hieu Nguyen, David, Soler Bartomeu, Tony CY Pang, Adam Wecker, Yifan Xiong, Fanfei Li, Lukas S. Huber, Joshua Jaeger, Romano De Maddalena, Xing Han Lù, Yuhui Zhang, Claas Beger, Patrick Tser Jern Kon, Sean Li, Vivek Sanker, Ming Yin, Yihao Liang, Xinlu Zhang, Ankit Agrawal, Li S. Yifei, Zechen Zhang, Mu Cai, Yasin Sonmez, Costin Cozianu, Changhao Li, Alex Slen, Shoubin Yu, Hyun Kyu Park, Gabriele Sarti, Marcin Briański, Alessandro Stolfo, Truong An Nguyen, Mike Zhang, Yotam Perlitz, Jose Hernandez-Orallo, Runjia Li, Amin Shabani, Felix Juefei-Xu, Shikhar Dhingra, Orr Zohar, My Chiffon Nguyen, Alexander Pondaven, Abdurrahim Yilmaz, Xuandong Zhao, Chuanyang Jin, Muyan Jiang, Stefan Todoran, Xinyao Han, Jules Kreuer, Brian Rabern, Anna Plassart, Martino Maggetti, Luther Yap, Robert Geirhos, Jonathon Kean, Dingsu Wang, Sina Mollaei, Chenkai Sun, Yifan Yin, Shiqi Wang, Rui Li, Yaowen Chang, Anjiang Wei, Alice Bizeul, Xiaohan Wang, Alexandre Oliveira Arrais, Kushin Mukherjee, Jorge Chamorro-Padial, Jiachen Liu, Xingyu Qu, Junyi Guan, Adam Bouyamourn, Shuyu Wu, Martyna Plomecka, Junda Chen, Mengze Tang, Jiaqi Deng, Shreyas Subramanian, Haocheng Xi, Haoxuan Chen, Weizhi Zhang, Yinuo Ren, Haoqin Tu, Sejong Kim, Yushun Chen, Sara Vera Marjanović, Junwoo Ha, Grzegorz Luczyna, Jeff J. Ma, Zewen Shen, Dawn Song, Cedegao E. Zhang, Zhun Wang, Gaël Gendron, Yunze Xiao, Leo Smucker, Erica Weng, Kwok Hao Lee, Zhe Ye, Stefano Ermon, Ignacio D. Lopez-Miguel, Theo Knights, Anthony Gitter, Namkyu Park, Boyi Wei, Hongzheng Chen, Kunal Pai, Ahmed Elkhanany, Han Lin, Philipp D. Siedler, Jichao Fang, Ritwik Mishra, Károly Zsolnai-Fehér, Xilin Jiang, Shadab Khan, Jun Yuan, Rishab Kumar Jain, Xi Lin, Mike Peterson, Zhe Wang, Aditya Malusare, Maosen Tang, Isha Gupta, Ivan Fosin, Timothy Kang, Barbara Dworakowska, Kazuki Matsumoto, Guangyao Zheng, Gerben Sewuster, Jorge Pretel Villanueva, Ivan Rannev, Igor Chernyavsky, Jiale Chen, Deepayan Banik, Ben Racz, Wenchao Dong, Jianxin Wang, Laila Bashmal, Duarte V. Gonçalves, Wei Hu, Kaushik Bar, Ondrej Bohdal, Atharv Singh Patlan, Shehzaad Dhuliawala, Caroline Geirhos, Julien Wist, Yuval Kansal, Bingsen Chen, Kutay Tire, Atak Talay Yücel, Brandon Christof, Veerupaksh Singla, Zijian Song, Sanxing Chen, Jiaxin Ge, Kaustubh Ponkshe, Isaac Park, Tianneng Shi, Martin Q. Ma, Joshua Mak, Sherwin Lai, Antoine Moulin, Zhuo Cheng, Zhanda Zhu, Ziyi Zhang, Vaidehi Patil, Ketan Jha, Qiutong Men, Jiaxuan Wu, Tianchi Zhang, Bruno Hebling Vieira, Alham Fikri Aji, Jae-Won Chung, Mohammed Mahfoud, Ha Thi Hoang, Marc Sperzel, Wei Hao, Kristof Meding, Sihan Xu, Vassilis Kostakos, Davide Manini, Yueying Liu, Christopher Toukmaji, Jay Paek, Eunmi Yu, Arif Engin Demircali, Zhiyi Sun, Ivan Dewerpe, Hongsen Qin, Roman Pflugfelder, James Bailey, Johnathan Morris, Ville Heilala, Sybille Rosset, Zishun Yu, Peter E. Chen, Woongyeong Yeo, Eeshaan Jain, Ryan Yang, Sreekar Chigurupati, Julia Chernyavsky, Sai Prajwal Reddy, Subhashini Venugopalan, Hunar Batra, Core Francisco Park, Hieu Tran, Guilherme Maximiano, Genghan Zhang, Yizhuo Liang, Hu Shiyu, Rongwu Xu, Rui Pan, Siddharth Suresh, Ziqi Liu, Samaksh Gulati, Songyang Zhang, Peter Turchin, Christopher W. Bartlett, Christopher R. Scotese, Phuong M. Cao, Aakaash Nattanmai, Gordon McKellips, Anish Cheraku, Asim Suhail, Ethan Luo, Marvin Deng, Jason Luo, Ashley Zhang, Kavin Jindel, Jay Paek, Kasper Halevy, Allen Baranov, Michael Liu, Advaith Avadhanam, David Zhang, Vincent Cheng, Brad Ma, Evan Fu, Liam Do, Joshua Lass, Hubert Yang, Surya Sunkari, Vishruth Bharath, Violet Ai, James Leung, Rishit Agrawal, Alan Zhou, Kevin Chen, Tejas Kalpathi, Ziqi Xu, Gavin Wang, Tyler Xiao, Erik Maung, Sam Lee, Ryan Yang, Roy Yue, Ben Zhao, Julia Yoon, Sunny Sun, Aryan Singh, Ethan Luo, Clark Peng, Tyler Osbey, Taozhi Wang, Daryl Echeazu, Hubert Yang, Timothy Wu, Spandan Patel, Vidhi Kulkarni, Vijaykaarti Sundarapandiyan, Ashley Zhang, Andrew Le, Zafir Nasim, Srikar Yalam, Ritesh Kasamsetty, Soham Samal, Hubert Yang, David Sun, Nihar Shah, Abhijeet Saha, Alex Zhang, Leon Nguyen, Laasya Nagumalli, Kaixin Wang, Alan Zhou, Aidan Wu, Jason Luo, Anwith Telluri, Summer Yue, Alexandr Wang, and Dan Hendrycks. Humanity's last exam, 2025. URL https://arxiv.org/abs/2501.14249.

Zile Qiao, Guoxin Chen, Xuanzhong Chen, Donglei Yu, Wenbiao Yin, Xinyu Wang, Zhen Zhang, Baixuan Li, Huifeng Yin, Kuan Li, Rui Min, Minpeng Liao, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. Webresearcher: Unleashing unbounded reasoning capability in long-horizon agents, 2025. URL https://arxiv.org/abs/2509.13309.

Hannah Rashkin, Vitaly Nikolaev, Matthew Lamm, Lora Aroyo, Michael Collins, Dipanjan Das, Slav Petrov, Gaurav Singh Tomar, Iulia Turc, and David Reitter. Measuring Attribution in Natural

- Language Generation Models. *Computational Linguistics*, pp. 1–64, 08 2023. ISSN 0891-2017. doi: 10.1162/coli_a_00486. URL https://doi.org/10.1162/coli_a_00486.
 - Aymeric Roucher, Albert Villanova del Moral, Merve Noyan, Thomas Wolf, and Clémentine Fourrier. Open-source DeepResearch Freeing our search agents, January 2025. URL https://huggingface.co/blog/open-deep-research.
 - Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. REPLUG: Retrieval-augmented black-box language models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 8371–8384, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.463. URL https://aclanthology.org/2024.naacl-long.463/.
 - Shuang Sun, Huatong Song, Yuhao Wang, Ruiyang Ren, Jinhao Jiang, Junjie Zhang, Fei Bai, Jia Deng, Wayne Xin Zhao, Zheng Liu, Lei Fang, Zhongyuan Wang, and Ji-Rong Wen. Simpledeepsearcher: Deep information seeking via web-powered reasoning trajectory synthesis, 2025. URL https://arxiv.org/abs/2505.16834.
 - Zhengwei Tao, Jialong Wu, Wenbiao Yin, Junkai Zhang, Baixuan Li, Haiyang Shen, Kuan Li, Liwen Zhang, Xinyu Wang, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. Webshaper: Agentically data synthesizing via information-seeking formalization, 2025. URL https://arxiv.org/abs/2507.15061.
 - Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. Executable code actions elicit better llm agents. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024.
 - 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, 2025. URL https://arxiv.org/abs/2504.12516.
 - Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He, Deyu Zhou, Pengjun Xie, and Fei Huang. WebWalker: Benchmarking LLMs in web traversal. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pp. 10290–10305, Vienna, Austria, July 2025a. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.508. URL https://aclanthology.org/2025.acl-long.508/.
 - Xixi Wu, Kuan Li, Yida Zhao, Liwen Zhang, Litu Ou, Huifeng Yin, Zhongwang Zhang, Yong Jiang, Pengjun Xie, Fei Huang, Minhao Cheng, Shuai Wang, Hong Cheng, and Jingren Zhou. Resum: Unlocking long-horizon search intelligence via context summarization, 2025b. URL https://arxiv.org/abs/2509.13313.
 - xAI. Grok 3 beta the age of reasoning agents, February 2025. URL https://x.ai/news/grok-3.
 - Ziyi Xia, Kun Luo, Hongjin Qian, and Zheng Liu. Open data synthesis for deep research, 2025. URL https://arxiv.org/abs/2509.00375.
 - 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. URL https://arxiv.org/abs/2210.03629.
 - Howard Yen, Tianyu Gao, Minmin Hou, Ke Ding, Daniel Fleischer, Peter Izsak, Moshe Wasserblat, and Danqi Chen. Helmet: How to evaluate long-context language models effectively and thoroughly. In *International Conference on Learning Representations (ICLR)*, 2025.
 - Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments, 2025. URL https://arxiv.org/abs/2504.03160.

A APPENDIX

A.1 EXISTING FRAMEWORKS

REACT (Yao et al., 2023) is a simple framework that allows an LLM agent to alternate between thinking and acting. This framework allows the agent to use tool calls across many turns. Following the original work's knowledge-intensive task settings, our implementation gives the LLM access to a single search tool—given a query, the tool returns a list of top 10 search results, from a search engine, along with their web contents. The search results are then concatenated and appended to the agent context for subsequent steps. When the LLM chooses not to use the search tool, the final output is used for evaluation.

In our implementation, we vary the maximum number of turns in each trajectory from 1 to 10. Consistently with SLIM, we use Google as the search engine, accessed through the Serper API, which returns a list of top 10 search results. Each search result contains a title, an URL, and a short snippet of the content. After obtaining the top 10 search results, we emulate previous RAG approaches by scraping all search result URLs and concatenate their content. Similar to SLIM, we use crawl4ai to scrape web pages. We truncate each scraped document to at most 10,000 characters, which corresponds to roughly 1,000 tokens.

We notice that REACT often hit the context window limit as the retrieval results are often too long. When the LLM API call fails due to the context window limit, we fallback to not using any tools and just ask the base LLM to answer the question. As a result, we only experiment with up to 10 turns, where the framework already falls back to not using any tools for most queries. A sketch of the framework is shown in Alg. 1.

SEARCH-01 (Li et al., 2025b) builds upon REACT with an additional "reason-in-document" step, where an LLM summarizes the list of top 10 search results and their contents before appending the results to the agent's input context. Although the summary added to the agent context is relatively short compared to the full search result, this approach still uses a large amount of browse calls in each search step, and the summarization steps incur a large amount of LLM token usage. In our setting, we vary the maximum number of turns in each trajectory from 1 to up to 100 turns.

Similar to REACT, the retrieval tool at each step is consisted of a single Serper API call, followed by multiple scraping operations. We adopt the code from the original implementation⁴, which uses BeautifulSoup⁵ to scrape the search result URLs. In this implementation, the scraping operation will extract part of the web content that best matches the short snippet returned by the search engine. The matching is done by simply computing the F1 scores between the snippet and sentences in the web page. Subsequently, the context is filled up with at most 2,500 characters from the web page. Then, all context from the search results are concatenated and appended to the agent context for the summarization step.

It's important to note that the scraping operation is relatively expensive due to the network latency, resulting in long running time for the framework. A sketch of the framework is shown in Alg. 2.

HuggingFace OpenDeepResearch (HF-ODR; Roucher et al., 2025) leverages a hierarchical structure consisting of a manager agent and a search agent. The manager agent calls the search agent to perform detailed searches, and the search agent iteratively interacts with the search engine and a simulated browser to gather information. When the search agent concludes its searches, it generates a summary of its searches and returns it to the manager agent. The manager agent may use the summary to issue additional queries or output the final answer. Furthermore, another key feature of HF-ODR is its access to additional tools, such as a Python interpreter. We use the default settings⁶, which fixes the maximum number of turns for the manager and search agent to be 20. A sketch of the framework is shown in Alg. 3. Specific descriptions of each tool can be found in Appendix A.2.

⁴https://github.com/RUC-NLPIR/Search-o1

⁵https://beautiful-soup-4.readthedocs.io/en/latest/

 $^{^6 \}verb|https://github.com/huggingface/smolagents/tree/main/examples/open_deep_research$

GPT-Researcher (GPT-R; Elovic, 2023) is a complex multi-agent system where each agent has distinct roles. Specifically, the system consists of a researcher conductor that orchestrates the search process, a report generator that generates the final report at the end of the search process, a context manager that summarizes search results, and a source curator that selects relevant sources from scraped web pages. Finally, GPT-R uses a deep researcher agent that acts as the node of a search tree, where each node is able to spawn multiple children nodes, each of which is a system with the previously described components. We use the default settings of the framework⁷, which fixes the number of depths of the search tree to be 2 and the breath of search at each depth to be 4. A sketch of the framework is shown in Alg. 4.

Algorithm 1: ReAct

864

865

866

867

868

870

871

872 873

874 875

876

877 878

879

880

882

883

884

885

887

889

890

891

892

893

894 895

896

897

899

900 901

902

903

904 905

906

907

908

909

910 911

912

913

914 915

916

917

```
Data: Task input x, LLM \theta, maximum number of turns T
Function search(q):
    return (title<sub>i</sub>, url<sub>i</sub>, snippet<sub>i</sub>)_{i=1}^k;
Function visit(u, q):
     D \leftarrow \operatorname{scrape}(u);
     return D[: 10000];
Result: Task output y
Turn t \leftarrow 1;
Context C \leftarrow \{x\};
\mathcal{T} \leftarrow \{\text{search}\};
while t < T do
     o_t \leftarrow \theta(C; \mathcal{T});
                                                /\star LLM may only call the search tool \star/
     switch o_t do
         case search do
               R \leftarrow \operatorname{search}(o_t);
                                                                                     /* Perform search */
               C \leftarrow C \cup \{o_t\};
                                          /* Visit every search result and append */
              for (t_i, u_i, s_i) \in R do
                   C \leftarrow C \cup \operatorname{visit}(u_i, s_i)
          case Final Answer do
              return o_t;
    t \leftarrow t + 1;
return \theta(C; final \ answer);
```

A.2 HUGGINGFACE OPEN DEEP RESEARCH TOOLS

HF-ODR is a hierarchical framework that consists of a manager agent and a search agent. The manager agent has access to the following tools:

- 1. **Search Agent**: an agent that will search the internet to answer a question.
- 2. **Visualizer**: given the path to a downloaded image, it will call and LLM to answer questions about the image.
- 3. **Text Inspector**: given the path to a downloaded text file, it will call and LLM to answer questions about the text.

The search agent has access to the following tools:

- Google Search: a search engine that will search the internet to answer a question. This tool
 uses Serper API in the backend.
- 2. Visit Tool: visit a URL and render the page in HTML as in a browser.
- 3. Page Up: navigate the current page by scrolling up.
- 4. **Page Down**: navigate the current page by scrolling down.

https://github.com/assafelovic/gpt-researcher

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

940

941

942

943

944

945 946

947 948 949

951 952

953

954

955 956

957 958

959 960

961

962

963

964 965

966

967

968

969

970

971

Algorithm 2: Search-o1 **Data:** Task input x, LLM θ , maximum number of turns T, summary interval n **Function** search(q): **return** (title_i, url_i, snippet_i) $_{i=1}^k$; **Function** visit(u, q): $D \leftarrow \operatorname{scrape}(u);$ $D \leftarrow \operatorname{split}(D) = \{d_i\}_{i=1}^m;$ if $q = \emptyset$ then **return** $d' \leftarrow d_1$; else $d' \leftarrow \arg\max_{d_i \in D} F1(d_i, q);$ return d'; **Result:** Task output y Turn $t \leftarrow 1$; Context $C \leftarrow \{x\}$; $\mathcal{T} \leftarrow \{\text{search}\};$ while t < T do $o_t \leftarrow \theta(C; \mathcal{T})$; /* LLM may only call the search tool */ switch o_t do case search do $R \leftarrow \operatorname{search}(o_t)$; /* Perform search */ $l \leftarrow \text{length}(C);$ $D \leftarrow \{c_i\}_{i=l-5}^l;$ for $(t_i, u_i, s_i) \in R$ do $D \leftarrow D \cup \operatorname{visit}(u_i, s_i)$; /* Visit every search result */ $C \leftarrow C \cup \{o_t, \theta(D; \text{summarize})\};$ case Final Answer do return o_t ; $t \leftarrow t + 1$; **return** $\theta(C; final \ answer);$

- 5. **Finder Tool**: find a text in the current page.
- 6. **Find Next**: find the next occurrence of the text in the current page.
- 7. **Archive Search**: search the archives for information.
- 8. **Text Inspector**: given the path to a downloaded text file, it will call and LLM to answer questions about the text.

Detailed descriptions of each tool can be found in the original implementation⁸.

A.3 TRAJECTORY-LEVEL ANALYSIS

In this subsection, we describe how we annotate each trajectory with the failure modes. For LLM-as-a-judge approaches, we use 03-2025-04-16 as the judge model. In each of the following LLM-as-a-judge approaches, we use the same judge model, and force the model to generate its response in a json format for easy parsing. We find that existing frontier LLMs are powerful enough to reliably check for simple yes/no questions and output them in a json format.

Confirmation bias. Confirmation bias occurs when the system finds a potential candidate that is incorrect in its search process, and subsequently spends the majority of its search budget on the same candidate without considering other options, leading to a lack of exploration in the search space. To this end, we first collect all the search queries that the system has made and then use an LLM to check if the search queries overly focus on a single wrong candidate. The judge model is given access

 $^{^8}$ https://github.com/huggingface/smolagents/blob/main/src/smolagents/default_tools.py

1015

1016

1017

1018

1019 1020

1021

1023

1024

1025

```
972
          Algorithm 3: HuggingFace Open Deep Research
973
          Data: Task input x, LLM \theta, maximum number of turns for search and main agents T_s and T_m,
974
                  respectively, and planning interval p
975
          web\_tools \leftarrow
976
           {Search, Visit, Page Up, Page Down, Finder, Find Next, Archive Search, Text Inspector};
977
          main\_tools \leftarrow \{search\_agent, Visualize, Text Inspector\};
978
          Function plan(q, c):
979
              /* Prompt the LLM to generate a plan
                                                                                                                       */
980
              return \theta(q, c; plan);
981
          Function search_agent(q):
982
              P \leftarrow \mathsf{plan}(q, \emptyset);
983
              C \leftarrow \{q, P\};
984
              t \leftarrow 1;
985
              while t < T_s do
                   if t \mod p = 0 then
986
                       P \leftarrow \mathsf{plan}(q, C);
987
                      C \leftarrow C \cup \{P\};
988
989
                   o_t \leftarrow \theta(C; \text{web\_tools});
                   if type(o_t) = final\_answer then
990
                      return o_t;
991
992
                   /* do the tool call, see A.2 for tool details
                                                                                                                       */
993
                   C \leftarrow C \cup \{o_t, \mathsf{tool}(o_t)\};
                   t \leftarrow t + 1;
994
995
              return \theta(C; final \ answer);
996
          Result: Task output y
997
          Turn t \leftarrow 1;
998
          P \leftarrow \mathsf{plan}(x,\emptyset);
999
          Context C \leftarrow \{x, P\};
1000
          /* the main agent plans and calls the search agent
                                                                                                                       */
          while t < T_m do
1001
              if t \mod p = 0 then
1002
                  P \leftarrow \mathsf{plan}(x,C);
1003
                  C \leftarrow C \cup \{P\};
1004
              o_t \leftarrow \theta(C; \text{main\_tools});
1005
              if type(o_t) = final\_answer then
                return o_t;
1007
              C \leftarrow C \cup \{o_t, \mathsf{tool}(o_t)\};
1008
              t \leftarrow t + 1;
1010
          return \theta(C; final \ answer);
```

to the groundtruth answer and the search queries, so it's able to determine if the search queries are relevant to the groundtruth answer and the similarities between different search queries. We consider a trajectory to have confirmation bias if a majority of the search queries are similar to each other, and focuses on a single wrong candidate. The prompt used for confirmation bias detection is shown in Table 5.

Unfocused search. Unfocused search occurs when the system generates overly generic search queries that are not useful for narrowing down the search space—the system cannot make any progress towards finding useful information. To this end, we first collect all the search queries that the system has made and then use an LLM to check if the search queries are generic and not useful for narrowing down the search space. We consider a trajectory to have unfocused search if a majority of the search queries are overly generic and not useful for narrowing down the search space. The prompt used for unfocused search detection is shown in Table 6.

```
1026
          Algorithm 4: GPT-Researcher
1027
          Data: Task input x, LLM \theta, research depth D, research breadth B, summary interval n
1028
          Function search(q):
1029
              return (title<sub>i</sub>, url<sub>i</sub>, snippet<sub>i</sub>)_{i=1}^k;
1030
          Function visit(u, q):
1031
               D \leftarrow \operatorname{scrape}(u);
1032
               D \leftarrow \operatorname{split}(D) = \{d_i\}_{i=1}^m;
1033
               if q = \emptyset then return d' \leftarrow d_1;
1034
               else d' \leftarrow \arg\max_{d_i \in D} F1(d_i, q);
1035
               return d';
1036
          Function plan(q):
1037
               /* Prompt the LLM to generate a list of gueries
                                                                                                                          */
               R \leftarrow \operatorname{search}(q);
1039
               return \theta(x, R; plan);
1040
          Function conduct_research(q):
1041
               /* Conduct research on one query by generating subqueries and
1042
                    retrieve and scrape
1043
               Q \leftarrow \mathsf{plan}(q);
               R \leftarrow \emptyset:
1045
               for q_i \in Q do
1046
                   for t_i, u_i, s_i \in search(q_i) do
1047
                       r_i \leftarrow \text{visit}(u_i, s_i);
1048
                       R \leftarrow R \cup r_i;
1049
               return \theta(x, R; process);
1050
          Function deep\_research(q, d):
1051
               /* Recursively plan and conduct research
                                                                                                                          */
1052
               Q \leftarrow \mathsf{plan}(q);
1053
               R \leftarrow \emptyset:
1054
               for q_i \in Q do
                   r_i \leftarrow \text{conduct\_research}(q_i);
1056
                    /* Prompt the LLM to generate takeaways and follow up
1057
                         questions
1058
                   q_i' \leftarrow \theta(r_i; \text{process});
                   if d < D then
                       R \leftarrow R \cup \text{deep\_research}(q'_i, d+1);
              return R;
1062
          Result: Task output y
1063
          Turn t \leftarrow 1;
1064
          Context C \leftarrow \{x\};
          P \leftarrow \mathsf{plan}(x);
          R \leftarrow \text{deep\_research}(P, 1);
1067
          return \theta(R; write report);
1068
```

Inefficient tool usage. Inefficient tool usage occurs when the system does not discover new information with its tool calls, and is therefore wasting its tool budget. Specifically, we use URLs as a proxy for the information discovered by the system—a tool call that only return URLs seen in previous search results is considered as a waste of tool budget. We use a simple heuristic for this analysis—iterate over all search calls made in the trajectory and keep track of seen URLs. Then, we report the percentage of search calls that only return URLs seen in previous search results.

1069 1070 1071

1072

1074

1075

1077

1078

1079

Groundtruth ignored. Groundtruth ignored occurs when the system encounters the correct answer in its search process, but does not use it to answer the question. One possible explanation is that the system is distracted by other noisy information in its context, preventing it from correctly identifying

Prompt for Confirmation Bias Detection

You are a helpful assistant that can analyze the trajectory of an information-seeking agent. You are given the search history of an agent and the correct answer for the question given to the agent. You should analyze the search history and determine if the agent spends more than half of the tool calls searching for the same incorrect answer. That is, the agent continues searching for the same topic even though it's not the correct answer, and spends half or more of its tool calls on these searches. Output your final conclusion with your reasoning and a single word: 'yes' if the agent spends more than half of its tool calls on the same incorrect answer or 'no' if the agent does not.

Reasoning: explain what the agent did, and if it did or did not focus its searches on a wrong answer.

Conclusion: "yes" or "no".

Search queries: <search-queries>
The correct answer is: <correct-answer>

Table 5: System prompt used for detecting confirmation bias in agent trajectories

Prompt for Unfocused Search Detection

You are a helpful assistant that can analyze the trajectory of an information-seeking agent. You are given the search history of an agent and the correct answer for the question given to the agent. You should analyze the search history and determine if the search queries are too broad. That is, the agent spends more than half of its tool calls on searches that are too generic and not specific enough to the question, and does not use enough tool calls to properly narrow down the search space by either eliminating wrong answers or verifying the correct answer. Output your final conclusion with your reasoning and a single word: 'yes' if the searches are too broad or 'no' if the searches are focused enough.

Reasoning: explain what the agent did, and if it did or did not use tool calls to properly narrow down the search space.

Conclusion: "yes" or "no".

Search queries: <search-queries>

Table 6: System prompt used for detecting unfocused search in agent trajectories

the groundtruth. We employ a simple approach for this analysis—we check if the groundtruth answer is present in any of the tool responses. We employ a LLM judge to enable fuzzy matching between the groundtruth answer and the tool responses. The prompt used for groundtruth ignored detection is shown in Table 7. We iterate over all tool calls and use this check to determine if any tool responses contain the groundtruth answer. We terminate the iteration if we find a tool response that contains the groundtruth answer, and report the percentage trajectories where at least one tool response contains the groundtruth answer.

Giving up. Giving up occurs when the system does not attempt to answer the question due to the lack of information in its context. Existing LLMs can often refuse to answer the question if it is not confident in answering the question, but this behavior is not desirable for information-seeking agents that could leverage additional tool calls to find the necessary information. We use a simple LLM judge to check if the system attempted to answer the question. The prompt used for giving up detection is shown in Table 8.

Hallucination. Hallucination occurs when the system generates information that is not supported by the information it has discovered in its search process. In deep research systems, it is not desirable to hallucinate information, as it could result in incorrect and misleading answers and thus affect the trustworthiness of the system. Inspired by previous works(Rashkin et al., 2023; Bohnet et al., 2022; Gao et al., 2023), we check if the system hallucinates information by first decomposing the model's explanation into a set of atomic claims. Then, we iterate through all the tool responses from the search process and check if the tool responses support all the claims. As long as one tool response support a claim, we consider the system to not have hallucinated that claim. In the end, we report the average percentage of unsupported claims in across trajectories. The prompt used for decomposing the model's explanation into a set of atomic claims is shown in Table 9, and the prompt

Prompt for Groundtruth Ignored Detection

You are a helpful assistant that can analyze the trajectory of an information-seeking agent. You are given a list of webpages and the correct answer for the question given to the agent. You should analyze the web contents and determine if it contains the correct answer. The correct answer is considered to be found if there are some context in the search results that is either a direct or near-exact match to the correct answer. Output your final conclusion with your reasoning and a single word: 'yes' if the content contains the correct answer or 'no' if the content does not contain the correct answer.

Reasoning: explain if the web content contains the correct answer.

Conclusion: "yes" or "no".
<tool-responses>

The correct answer is: <correct-answer>

Table 7: System prompt used for detecting groundtruth ignored in agent trajectories

Prompt for Giving Up Detection

You are a helpful assistant that can analyze the final output of an information-seeking agent. You are to check if the agent decides that it cannot find the correct answer. For example, if the explanation states that it cannot find enough relevant information to answer the question, or if the response is simply empty or "I don't know", then the agent did not attempt to answer the question. Output your final conclusion with a single word "yes" if the agent decides it did not find enough information to answer the question or "no" otherwise.

Conclusion: "yes" or "no".
Final output: <final-output>

Table 8: System prompt used for detecting giving up in agent trajectories

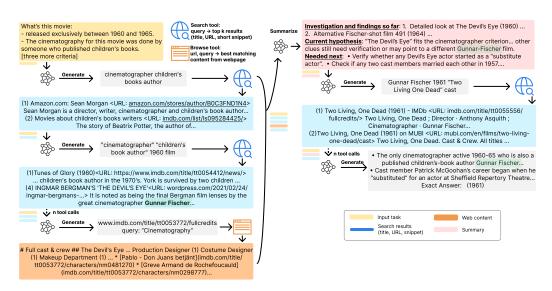


Figure 8: An example of a SLIM trajectory.

used for hallucination detection is shown in Table 10. These prompts are derived from previous works that show LLMs can reliably decompose texts into a set of atomic claims and check if claims are supported by a piece of text—they also achieve high agreement with human judges (Gao et al., 2023; Kamoi et al., 2023; Yen et al., 2025).

A.4 SLIM DETAILS

We show an example of a SLIM trajectory in Figure 8. Furthermore, a sketch of the framework is shown in Alg. 5.

1188 **Prompt for Decomposing Explanation into Atomic Claims** 1189 Read the given explanation and generate a list of atomic claims that are supported by 1190 the explanation. Atomic claims that are basic facts that cannot be further broken down. 1191 Generate at most 10 claims for the explanation. Use the following as an example: 1192 1193 Explanation: Searching UFCStats for featherweight bouts 1194 where the loser landed 14 of 83 significant strikes (16.87 %) 1195 and went 0-for-4 on takedowns returns the fight Myles Jury vs. Ricky Glenn at UFC 219: Cyborg vs Holm (30 Dec 2017). 1196 • Ricky Glenn (nickname "The Gladiator" | a synonym 1197 for swordsman) 1198 was the loser: sig. strikes 14/83 (16.87 %), takedowns 0/4. 1199 • Both fighters (Jury 29, Glenn 28) were under 35 and are American. 1201 · The referee was John McCarthy, whose first event for the UFC was in 1994. 1203 Thus, the MMA event is UFC 219: Cyborg vs Holm. 1205 Exact Answer: UFC 219: Cyborg vs Holm 1207 Confidence: 75% 1208 **Atomic Claims:** 1209 - Ricky Glenn was the loser - Ricky Glenn was nicknamed "The Gladiator" 1210 - The sig. strike rate of Ricky Glenn was 14/83 (16.87- The takedown rate of Ricky Glenn 1211 was 0/4 1212 - Jury was age 29 1213 - Glenn was age 28 1214 - Jury is American 1215 - Glenn is American 1216 - The referee was John McCarthy 1217 - John McCarthy's first event for the UFC was in 1994 1218 Output the atomic claims in the form of a json list.

Table 9: System prompt used for decomposing the model's explanation into a set of atomic claims

Prompt for Hallucination Detection

You are a helpful assistant that can analyze the trajectory of an information-seeking agent. You are given a list of webpages and a list of claims made by the agent. You should analyze the web contents to determine if each claim is supported by the web content. A claim is supported by the web content if its factual information is mostly supported by the web content, and is not contradicted by the web content. Output your final conclusion with a list of claims that are supported by the web content. Output the list in the form of a json list, and you only need to write the index of the supported claims in the list and nothing else.

Webpages: <webpages>

1219 1220

1223

1224

1225

1226

1227

1228

1229

1230

1231 1232

1233

1236

1237

1239

1240

1241

Atomic Claims: <atomic-claims>

Table 10: System prompt used for detecting hallucination in agent trajectories

A.5 EXPERIMENTAL DETAILS

We use o3, o4-mini, and Claude-4-Sonnet as our base models. To calculate the cost, we use the prices listed in Table 11, which are obtained from respective websites https://platform.openai.com/docs/models/o3, https://platform.openai.com/docs/models/o4-mini, https://claude.com/pricing#api, https://www.firecrawl.dev/pricing.

To calculate the token cost, we take a weight sum of the token usage across all LLM calls: non-cached input tokens plus 4 times the total output tokens, and multiply the results by price per token. We

```
1242
           Algorithm 5: SLIM
1243
           Data: Task input x, LLM \theta, maximum number of turns T, summary interval n
1244
           Function search(q):
1245
               return (title<sub>i</sub>, url<sub>i</sub>, snippet<sub>i</sub>)_{i=1}^k;
1246
           Function visit(u, q):
1247
                D \leftarrow \text{scrape}(u);
1248
                 D \leftarrow \operatorname{split}(D) = \{d_i\}_{i=1}^m;
1249
                if q = \emptyset then return d' \leftarrow d_1;
1250
                else d' \leftarrow \arg \max_{d_i \in D} \text{ROUGE-L}(d_i, q);
1251
                return d';
1252
           Result: Task output y
1253
           Turn t \leftarrow 1;
1254
           Context C \leftarrow \{x\};
1255
           \mathcal{T} \leftarrow \{\text{search}, \text{visit}\};
1256
           while t < T do
1257
                if t \mod n = 0 then
1258
                  C \leftarrow \theta(C; \text{summarize});
                                                                                   /\star Summarize every n turns \star/
1259
                o_t \leftarrow \theta(C; \mathcal{T});
1260
                switch o_t do
1261
                      case search do
1262
                           q_t \leftarrow o_t;
1263
                           C \leftarrow C \cup \{o_t, \operatorname{search}(q_t)\};
1264
                      case visit do
1265
                           u_t, s_t \leftarrow o_t;
1266
                           C \leftarrow C \cup \{o_t, \operatorname{visit}(u_t, s_t)\};
1267
                      case Final Answer do
1268
                           return o_t;
1269
                t \leftarrow t + 1;
1270
1271
           return \theta(C; final \ answer);
```

exclude cached tokens from the calculation because in practice, long-horizon systems are expected to have a large amount of cached tokens and system implementation that takes advantage of caching. Then, for the total cost, we add in the number of search API and scrape URL operations, multiplied by their respective prices.

Table 11: Pricing for different components. Numbers are obtained from respective websites:

	Cost
03	\$2.0 / M token
o4-mini	\$1.1 / M token
Claude-4-Sonnet	\$3.0 / M token
Google search	\$0.5 / K query
Scrape URL	\$0.83 / K query

A.6 ADDITIONAL RESULTS

127212731274

1275

1276

1277

1278 1279

1289 1290

1291

Here we provide the concrete results for SLIM with different base models—o4-mini is shown in Table 12, and Claude-4-Sonnet is shown in Table 13.

Table 12: Main results with o4-mini as the base model. All results are macro-averaged across test instances. The number of tokens is shown in 10,000s. The cost is shown in US dollars. T denotes the maximum number of turns in each trajectory.

		BrowseComp				HLE				
	T	Score (†)	Tokens (↓)	Tools (↓)	Cost (↓)	Score (†)	Tokens (↓)	Tools (↓)	Cost (\1)	
o4	-	5.0	5.1	0.0	0.06	15.0	2.2	0.0	0.02	
	1	1.3	4.6	1.0	0.05	17.0	4.0	0.5	0.04	
REACT	5	3.0	7.7	2.1	0.09	15.3	4.6	0.7	0.05	
	10	2.3	7.4	2.3	0.08	15.3	4.9	0.8	0.05	
Search-o1	1	6.3	6.2	10.0	0.08	13.0	2.6	3.5	0.03	
	5	11.3	13.8	49.7	0.19	23.3	4.0	11.9	0.05	
	25	25.0	45.4	207.7	0.66	22.3	5.5	22.5	0.08	
	50	28.7	76.1	351.5	1.12	19.3	7.3	26.3	0.10	
	100	36.0	124.4	546.7	1.80	21.3	6.6	25.8	0.09	
HF-ODR	20	15.0	38.9	15.4	0.44	16.3	8.3	3.9	0.09	
GPT-R	-	4.0	8.5	82.5	0.16	11.3	9.7	100.8	0.19	
	10	14.0	5.7	8.8	0.07	21.0	3.6	3.1	0.04	
	25	24.3	24.0	23.2	0.28	23.7	7.2	5.9	0.08	
SLIM	50	31.0	73.7	40.1	0.83	25.7	10.0	7.0	0.1	
	100	34.0	92.9	45.2	1.05	26.7	12.2	7.7	0.14	
	150	37.0	107.8	49.5	1.22	24.7	14.4	8.6	0.10	

Table 13: Main results with Claude-4-Sonnet as the base model. All results are macro-averaged across test instances. The number of tokens is shown in 10,000s. The cost is shown in US dollars. T denotes the maximum number of turns in each trajectory.

			Browse	Comp		HLE			
	T	Score (†)	Tokens (↓)	Tools (↓)	Cost (↓)	Score (†)	Tokens (↓)	Tools (↓)	Cost (↓)
Claude-4-Sonnet	-	1.0	1.9	0.0	0.06	6.3	3.9	0.0	0.12
	1	0.3	0.0	0.0	0.00	8.3	0.0	0.0	0.00
REACT	5	0.3	0.0	0.0	0.00	8.3	0.0	0.0	0.00
	10	0.3	0.0	0.0	0.00	8.0	0.0	0.0	0.00
	1	2.0	1.5	9.0	0.05	10.0	2.9	10.0	0.09
	5	3.7	6.0	44.1	0.21	11.7	5.3	29.5	0.18
SEARCH-01	10	7.0	10.7	79.5	0.38	16.0	6.3	35.6	0.22
SEARCH-UI	25	8.0	20.1	149.9	0.72	13.0	6.8	41.1	0.24
	50	10.0	22.9	170.3	0.82	12.7	7.0	40.7	0.24
	100	10.0	19.4	148.3	0.70	12.3	6.4	38.5	0.22
HF-ODR	20	6.7	98.8	30.4	2.98	17.3	105.0	26.5	3.16
GPT-R	-	2.3	7.9	106.5	0.32	8.0	6.9	94.9	0.28
	10	2.7	2.8	8.9	0.09	10.3	2.5	6.9	0.08
	25	9.7	5.1	21.6	0.17	15.0	2.8	10.2	0.09
SLIM	50	10.0	5.0	27.1	0.16	17.3	3.0	9.9	0.10
	100	10.7	4.8	28.1	0.16	14.0	2.9	10.5	0.09
	150	10.0	5.2	30.7	0.17	16.7	3.1	11.1	0.10