

# AGENT-AS-TOOL: A STUDY ON THE HIERARCHICAL DECISION MAKING WITH REINFORCEMENT LEARNING

Yanfei Zhang

Independent Researcher

zyf909533@gmail.com

## ABSTRACT

Large Language Models (LLMs) have emerged as one of the most significant technological advancements in artificial intelligence in recent years. Their ability to understand, generate, and reason with natural language has transformed how we interact with AI systems. With the development of LLM-based agents and reinforcement-learning-based reasoning models, the study of applying reinforcement learning in agent frameworks has become a new research focus. However, all previous studies face the challenge of deciding the tool calling process and the reasoning process simultaneously, and the chain of reasoning was solely relied on the unprocessed raw result with redundant information and symbols unrelated to the task from the tool, which impose a heavy burden on the model’s capability to reason. Therefore, in our research, we proposed a hierarchical framework **Agent-as-tool** that detach the tool calling process and the reasoning process, which enables the model to focus on the verbally reasoning process while the tool calling process is handled by another agent. Our work had achieved comparable results with only a slight reinforcement fine-tuning on 180 samples, and had achieved exceptionally well performance in Bamboogle (Brown et al., 2020) with **63.2%** of exact match and **75.2%** in cover exact match, exceeding Search-R1 by 4.8% in exact match and 3.2% in cover exact match.

## 1 INTRODUCTION

Large Language Models (LLMs) have achieved remarkable progress in a wide range of natural language understanding and generation tasks (Liu et al., 2025; Zhang et al., 2024). As the complexity of tasks increases, a common approach is to augment LLMs with access to external tools, such as web search engines, calculators, or code interpreters. This tool-augmented paradigm enables agents to interact with the environment and perform planning, reasoning, and execution steps beyond the model’s pretraining distribution.

Recent advancements have explored integrating reinforcement learning (RL) into these agent frameworks, aiming to improve decision-making over tool usage and multi-hop reasoning steps (Guo et al., 2025; Jin et al., 2025). However, a major limitation remains: existing RL-enhanced agents conflate the tool invocation process with the verbal reasoning process. This tight coupling leads to several challenges: (1) The agent must learn tool selection, input construction, and reasoning jointly, which increases training difficulty and noise; (2) Reasoning often proceeds over noisy, unstructured outputs returned directly from external tools, which degrades answer quality.

To address these challenges, we propose **Agent-as-tool**, a hierarchical reasoning architecture in which reasoning and tool execution are explicitly decoupled as shown in Figure 1. The framework introduces a **Planner** and a **Toolcaller** as two separate agent components. The **Planner** focuses on natural language reasoning and high-level decision-making, while the **Toolcaller** is responsible for managing the tool interface (e.g., invoking web search) and returning structured observations.

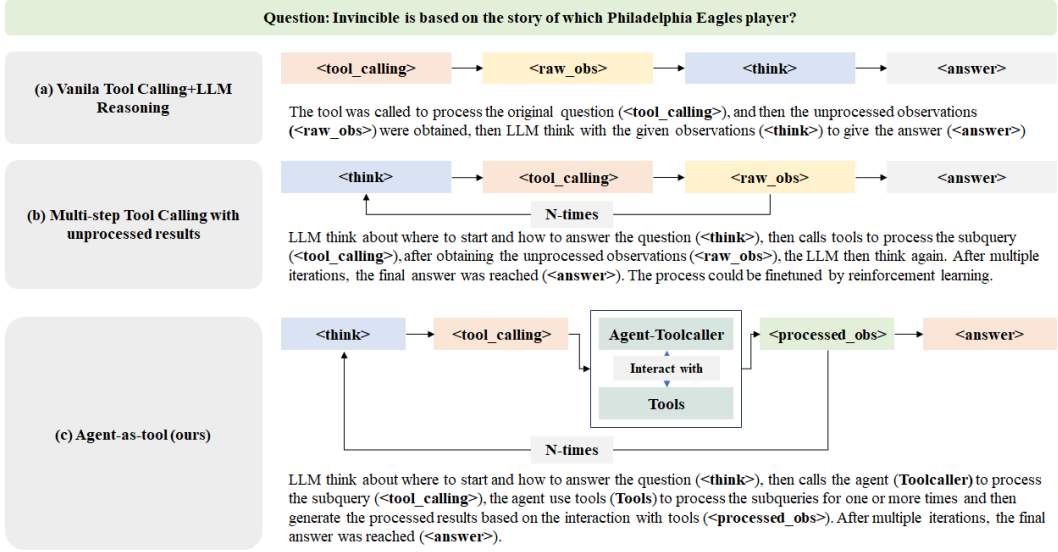


Figure 1: The trajectory of a single sample from a batch of questions processed in different research configurations. In our **Agent-as-tool** method, we employed the agent as a tool instead of calling the tool directly. The **Planner** is responsible for the tool calling process and the reasoning process, and the **Toolcaller** is responsible for the tool calling process to provide sufficient processed observations.

The advantages of this design are twofold: (1) It simplifies the RL optimization process by assigning each sub-agent a focused objective; (2) It improves reasoning accuracy by allowing the **Planner** to operate on cleaner, more structured inputs. Furthermore, we apply a lightweight reinforcement fine-tuning procedure using GRPO on just 180 samples to demonstrate the efficiency of our framework.

This paper makes the following contributions:

- We propose **Agent-as-tool**, a hierarchical agent framework that separates reasoning and tool usage via a **Planner** and a **Toolcaller**.
- We introduce a reinforcement learning protocol that enhances **Planner** behavior while masking **Toolcaller** outputs to preserve credit assignment integrity.
- We empirically validate our framework on multiple multi-hop QA datasets and achieve state-of-the-art performance on Bamboogle.
- We provide qualitative insights showing that hierarchical decoupling improves reasoning clarity and decomposition over existing baselines like Search-R1.

## 2 LITERATURE REVIEW

### 2.1 AGENT FRAMEWORKS BASED ON PRE-DEFINED REASONING STEPS

There are several agent researches that are designed to perform tasks with pre-trained LLMs and with pre-defined reasoning steps, including the CAMEL (Li et al., 2023a), OpenManus (Liang et al., 2025) and MetaGPT (Hong et al., 2023). These works tend to extend the capabilities of the pre-trained LLM with additional rule-based reasoning steps to ‘stimulate’ the internal reasoning capabilities of the LLM to achieve better performances.

Specifically, considering the search and information retrieval for task completion scenario, there are also considerable works, including the Search-o1 (Li et al., 2025a), OpenResearcher (Zheng et al., 2024) (majorly focusing on the scientific research scenario).

---

## 2.2 RL REASONING AGENTS

With the development of RL training frameworks and the Deepseek-R1 (Guo et al., 2025) setting, there are also considerable works to implement R1-style training paradigms on the LLM-based agents. The searching and information retrieval tasks were the first to be considered in this scenario, including the R1-searcher (Song et al., 2025) and DeepResearcher (Zheng et al., 2025).

There are also several works that integrate other external tools under the framework to complete different tasks, including the ToRL (Li et al., 2025b) that integrate the python interpreter tool, ToolRL (Qian et al., 2025) that flexibly integrate different toolkits with different pre-defined datasets (e.g. API-Bank (Li et al., 2023b)), SWiRL (Goldie et al., 2025) that control the tool selection process with different labels (`<calculator>` for calculator tool, and `<search_query>` for web search tool).

The generic process of calling an agent in these researches can be concluded as a sequence of thinking `<think>`, followed by a tool calling query enclosed with `<tool_query>`, then the tool returns observations `<obs>`. With the reasoning on each step, the final answer could be reached whenever the agent think the ground truths are sufficient enough to give the final answer. It is a much simpler configuration with a ReAct-like tool calling process (Yao et al., 2023), then reinforcement learning are applied to explore whether the model could exhibit the capabilities beyond simple reasoning to reach the next hop, as shown in Figure 1 as Multi-step Tool Calling with Unprocessed Results.

## 3 METHODOLOGY

We propose the **Agent-as-tool** framework as a hierarchical design for multi-hop reasoning tasks. It separates the planning and tool usage responsibilities between two agent components: a high-level **Planner** and a subordinate **Toolcaller**. The **Planner** manages reasoning and task decomposition, while the **Toolcaller** executes external actions such as web search. This section outlines the design of both components and the reinforcement learning procedure employed to optimize the **Planner**.

### 3.1 AGENT ARCHITECTURE

#### 3.1.1 PLANNER

The **Planner** is a language model agent responsible for high-level reasoning and tool invocation decisions. It reasons about the current task state and emits tool usage instructions in natural language.

**Reasoning:** The **Planner** conducts internal reasoning enclosed in `<think>...</think>` tags, in line with DeepSeek-R1 conventions (Guo et al., 2025). It uses previous observations and the original query to plan the next subtask.

**Tool Invocation:** Tool calls are expressed as sub-queries wrapped in `<tool_calling>...</tool_calling>` tags. These queries are interpreted by the **Toolcaller**, and the results are returned to the **Planner** as `<obs>...</obs>` blocks for further reasoning.

#### 3.1.2 TOOLCALLER

The **Toolcaller** is a dedicated LLM-based agent designed to interface with external tools. In our implementation, it wraps a web search tool and processes queries issued by the **Planner**.

We implement the **Toolcaller** using a CAMEL-style chat agent (Li et al., 2023a), powered by GPT-4o-mini (Hurst et al., 2024). It could retrieve top- $k$  search results multiple times and returns structured summaries to the **Planner**. Although our current prototype uses only web search, the architecture supports extension to tools like calculators or code interpreters, also including MCP-based tool servers.

---

### 3.2 REINFORCEMENT LEARNING WITH GRPO

#### 3.2.1 TRAINING OBJECTIVE

We employ Generalized Reinforcement Policy Optimization (GRPO) (Shao et al., 2024) to fine-tune the **Planner**. The objective is:

$$\mathcal{J}(\Theta) = \mathbb{E}_{x \sim \mathcal{D}, \{y_i\}_{i=1}^G \sim \pi_{\text{old}}(\cdot|x)} \left[ \frac{1}{G} \sum_{i=1}^G \left[ \min \left( \frac{\pi_{\Theta}(y_i|x)}{\pi_{\text{old}}(y_i|x)} A_i, \text{clip} \left( \frac{\pi_{\Theta}(y_i|x)}{\pi_{\text{old}}(y_i|x)}, 1 - \varepsilon, 1 + \varepsilon \right) A_i \right) - \beta D_{\text{KL}}(\pi_{\Theta} || \pi_{\text{ref}}) \right] \right] \quad (1)$$

where  $x$  is sampled from dataset  $\mathcal{D}$ ,  $y_i$  is a rollout,  $A_i$  is the advantage,  $\varepsilon$  is the clipping threshold, and  $\beta$  regulates KL penalty.

#### 3.2.2 OBSERVATION MASKING

To prevent reward leakage through Toolcaller-generated outputs, we mask the `<obs>` blocks during reward modeling and training. These segments are replaced with special token `<fim_pad>`, which is trained to embed close to zero.

#### 3.2.3 REWARD FUNCTION

Our reward function balances correctness and formatting constraints:

$$\text{Reward} = \begin{cases} \text{F1 score} & \text{if answer is correctly formatted} \\ -2 & \text{otherwise} \end{cases} \quad (2)$$

The model receives a high reward when generating a valid and correct response, and a penalty when output is malformed.

## 4 EXPERIMENTS

### 4.1 EXPERIMENT SETTINGS

#### 4.1.1 MODEL AND HYPERPARAMETERS

We use Qwen-2.5-7B-Instruct (Qwen et al., 2025) as our base model. The training is conducted by an customized implementation of rollout and a customized implementation of GRPO on trl (von Werra et al., 2020). At each training step, we sample a batch of training data from the training set and calculate the reward for each rollout. Then, we update the policy by maximizing the reward.

The batch size is set to 3 for each training step and each sample contains 12 rollouts for each prompt. Each rollout contains at most 10 rounds of tool calling.

#### 4.1.2 TRAINING SETTINGS

We conducted quite a small scale of training for the **Agent-as-tool**. We trained the **Agent-as-tool** for 60 steps with each step containing 3 training samples, and each training sample contains 12 rollouts, with total size of only 180 samples and 2160 rollouts. The training data entries were selected from the HotpotQA (Yang et al., 2018) and 2WikiMultiHopQA (Ho et al., 2020) datasets with the same ratio as the R1-searcher (Song et al., 2025).

During the training process, we observed that the loss of the **Agent-as-tool** is not stable for the first 30 steps, which is likely due to the small training data, but after 30 steps, the loss is stable and close to 0 and the performance of the **Agent-as-tool** also stabilized.

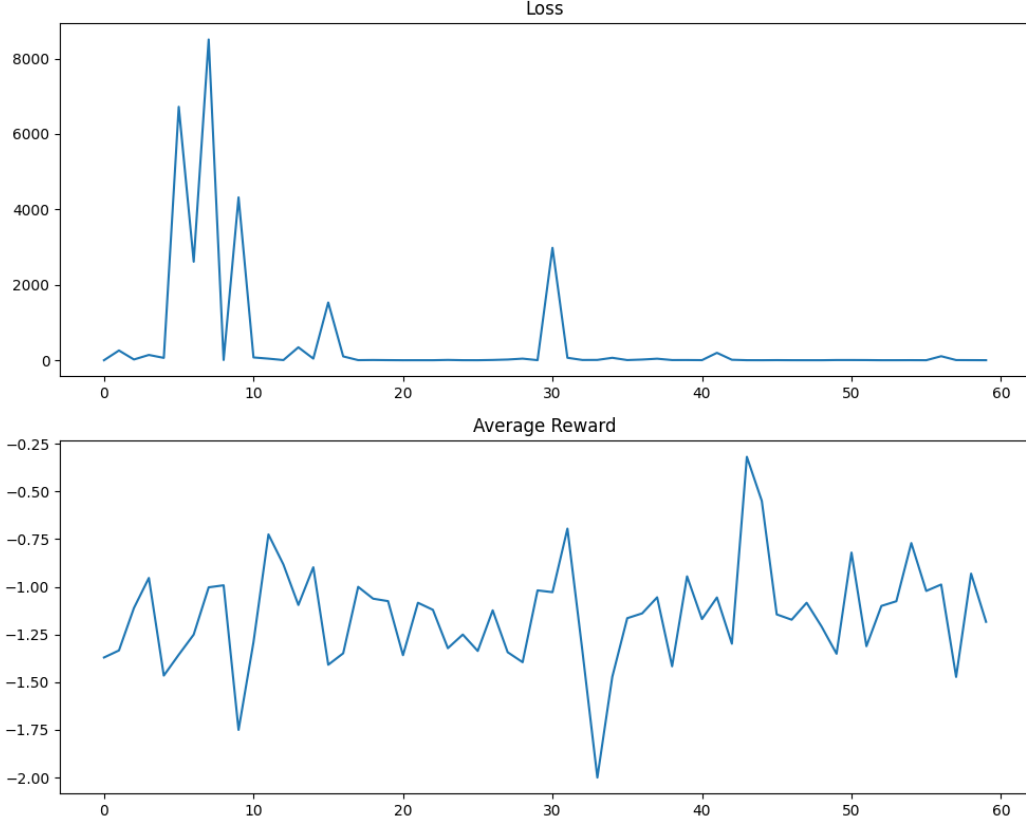


Figure 2: Training progress of the Agent-as-tool model showing loss convergence over training steps. The loss becomes stable after approximately 30 training steps.

The training curve illustrated in Figure 2 shows the convergence behavior of our model during the reinforcement learning process.

#### 4.1.3 BENCHMARK SETTINGS

In order to evaluate the performance of the **Agent-as-tool**, we conducted experiments on the open-domain question-answering task. We selected multiple multi-hop reasoning tasks to evaluate the performance of the **Agent-as-tool**, including the HotpotQA (Yang et al., 2018), 2WikiMultiHopQA (Ho et al., 2020), MuSiQue (Trivedi et al., 2022), and bamboogle (Press et al., 2023).

#### 4.1.4 BASELINE SETTINGS

We have 1 information retrieval tool: web search.

We then compare the performance of the **Agent-as-tool** with the following baselines:

- **direct IO**: This baseline employs the direct output of the Qwen-2.5-7B-Instruct (Qwen et al., 2025) as the answer without any external tool calling.
- **direct IO with web search**: This baseline employs the direct output of the Qwen-2.5-7B-Instruct (Qwen et al., 2025), but enables the web search to process the original question and return the top-k results as additional observations.
- **CAMEL Agent**: This baseline employs the CAMEL (Li et al., 2023a) chat agent driven by the GPT-4o-mini (Hurst et al., 2024).

- **CAMEL Agent with web search:** This baseline employs the CAMEL (Li et al., 2023a) chat agent driven by the GPT-4o-mini (Hurst et al., 2024) and the same tool setting as the **Agent-as-tool** with web search tool only. This baseline is used as the reference for multi-hop reasoning tasks conducted with the rule-based agent framework.
- **Search-R1:** We directly compare the performance of the **Agent-as-tool** with the Search-R1 (Jin et al., 2025) in our configurations with web search tool for a fair comparison. As **Search-R1** cannot be directly integrated with the CAMEL (Li et al., 2023a) chat agent, we would directly returns the search results as the answer instead of using another **Toolcaller**.

We conducted experiments with **Agent-as-tool** with pre-finetuned and post-finetuned models.

In align with the Deepseek-R1 setting (Guo et al., 2025), we adopted the same prompt setting for all the baselines and the **Agent-as-tool** except Search-R1 (Jin et al., 2025) that is equipped with its original prompt setting, and we also modified the tool calling process to enable Search-R1 to access the unprocessed web search results.

#### 4.1.5 EVALUATION METRICS

In this paper, we focus on the performance of the **Agent-as-tool** in terms of the correctness of the answer, therefore, we employed the exact match metric (**EM**), the cover exact match metric (**CEM**) to evaluate the performance of the **Agent-as-tool**.

## 4.2 QUANTITATIVE EXPERIMENT RESULTS

The qualitative results are shown in 1. Based on the results, we can see that the **Agent-as-tool** outperforms most of the baselines except for the EM metric in the HotpotQA, 2WikiMultiHopQA, and MuSiQue datasets, where Search-R1 still has the best performance. However, in terms of the CEM metric, our model has a substantial improvement over all the baselines, except in HotpotQA where Search-R1 still has the best performance (64.2% vs 57.4%). And in the Bamboogle dataset (Press et al., 2023), the **Agent-as-tool** with web search tool integrated to the **Toolcaller** (CAMEL (Li et al., 2023a) agent) achieves the best performance with EM of 63.2% and CEM of 75.2%.

Table 1: Performance Comparison Across Different Datasets

Dataset	Model	EM (%)	CEM (%)
Bamboogle	Direct IO	17.6	26.4
	Direct IO + Web Search	29.6	42.4
	CAMEL	36.8	47.2
	CAMEL + Web Search	51.2	62.4
	Search-R1 + Web Search	58.4	72.0
	<b>Agent-as-tool-Base + Web Search</b>	60.0	71.2
	<b>Agent-as-tool-Instruct + Web Search</b>	<b>63.2</b>	<b>75.2</b>
HotpotQA	Direct IO	20.0	27.2
	Direct IO + Web Search	32.6	52.8
	CAMEL	23.2	44.2
	CAMEL + Web Search	32.4	59.4
	Search-R1 + Web Search	<b>47.2</b>	<b>64.2</b>
	<b>Agent-as-tool-Base + Web Search</b>	35.0	55.2
	<b>Agent-as-tool-Instruct + Web Search</b>	37.2	57.4
2WikiMultiHopQA	Direct IO	22.6	25.4
	Direct IO + Web Search	27.2	40.2

*Continued on next page*

Table 1 – continued from previous page

Dataset	Model	EM (%)	CEM (%)
	CAMEL	20.8	34.6
	CAMEL + Web Search	35.0	69.4
	Search-R1 + Web Search	<b>52.4</b>	68.0
	<b>Agent-as-tool-Base + Web Search</b>	42.8	68.0
	<b>Agent-as-tool-Instruct + Web Search</b>	44.6	<b>70.0</b>
MuSiQue	Direct IO	4.8	9.0
	Direct IO + Web Search	14.0	18.0
	CAMEL	9.2	18.8
	CAMEL + Web Search	16.0	29.4
	Search-R1 + Web Search	<b>20.8</b>	28.6
	<b>Agent-as-tool-Base + Web Search</b>	15.6	28.8
	<b>Agent-as-tool-Instruct + Web Search</b>	18.4	<b>29.8</b>

We compared the performance of the **Agent-as-tool** before and after the reinforcement fine-tuning process. The table is shown in 2. Based on the results, we can see that the Reinforcement fine-tuning based on GRPO (Shao et al., 2024) substantially improves the performance of the **Agent-as-tool** in all datasets with an average improvement of 2.5% in EM and 2.3% in CEM.

Table 2: Performance improvements after reinforcement fine-tuning

Dataset	Pre-finetuned		Post-finetuned		Improvement	
	EM (%)	CEM (%)	EM (%)	CEM (%)	EM (%)	CEM (%)
Bamboogle	60.0	71.2	63.2	75.2	<b>+3.2</b>	<b>+4.0</b>
HotpotQA	35.0	55.2	37.2	57.4	<b>+2.2</b>	<b>+2.2</b>
2WikiMultiHopQA	42.8	68.0	44.6	70.0	<b>+1.8</b>	<b>+2.0</b>
MuSiQue	15.6	28.8	18.4	29.8	<b>+2.8</b>	<b>+1.0</b>
<b>Average</b>	38.4	55.8	40.9	58.1	<b>+2.5</b>	<b>+2.3</b>

Comparing with the CAMEL baseline with web search tool integrated (**CAMEL + Web Search**), the **Agent-as-tool** pre-finetuned and post-finetuned achieved a substantial improvement in EM and CEM, stating the necessity that the **Agent-as-tool** that enables the model to control when and what to be called in a tool calling is a more effective framework for multi-hop reasoning tasks.

Comparing with the Search-R1 baseline (**Search-R1 + Web Search**), which is the current best performing research of its kind, the **Agent-as-tool-Instruct** has substantial improvements over the Bamboogle dataset, which improves the EM by 4.8% and CEM by 3.2%, stating the effectiveness of the **Agent-as-tool** in multi-hop reasoning tasks. Besides, the **Agent-as-tool** conducted fine-tuning with 180 samples, which indicates the efficiency of the fine-tuning process.

### 4.3 QUALITATIVE RESULTS INSPECTION AND ANALYSIS

#### 4.3.1 COMPARISON OF THE AGENT-AS-TOOL AND THE SEARCH-R1

Comparing with Search-R1 baseline (**Search-R1 + Web Search**), the **Agent-as-tool-Instruct** had several advantages qualitatively:

- The **Agent-as-tool-Instruct** could reason with less fuzzy and more structured observations, comparing with the **Search-R1 + Web Search** which would need to reason with the unprocessed web search results with fuzzy and unstructured symbols or other unrelated details.

- As the **Agent-as-tool-Instruct** adopt a hierarchical reasoning process which segregate the reasoning process and the tool calling process, the agent could have a better linearly text-based reasoning process comparing with the **Search-R1 + Web Search**.

The qualitative comparison as a example is shown in Figure 3 (in Appendix). The **Agent-as-tool-Instruct** could reason with less fuzzy and more structured observations, comparing with the **Search-R1 + Web Search** which would need to reason with the unprocessed web search results with fuzzy and unstructured symbols or other unrelated details.

#### 4.3.2 COMPARISON OF THE RESULTS BEFORE AND AFTER THE REINFORCEMENT FINE-TUNING

Comparing with **Agent-as-tool-Base**, the **Agent-as-tool-Instruct** had several advantages qualitatively:

- The **Agent-as-tool-Instruct** identify the correct decomposition of the question to identify the first hop and the second hop to be solved by the agent, comparing with the **Agent-as-tool-Base** which would not be able to decompose the multi-hop question correctly so it directly feed the agent with the whole question (only a slightly change from the original manner). If the agent is not capable of reasoning the multi-hop question correctly, the **Agent-as-tool-Base** would not be able to answer the question correctly.
- As the **Agent-as-tool-Instruct** was instructed to reason with the agent powered by the pretrained model, the fine-tuned model could give a more structured and reasonable question to be answered by the agent comparing with the **Agent-as-tool-Base**.

The qualitative comparison as a example is shown in Figure 4 (in Appendix). The **Agent-as-tool-Instruct** could correctly decompose the question to identify the first hop and the second hop, comparing with the **Agent-as-tool-Base** which would not be able to decompose the multi-hop question correctly so it directly feed the agent with the whole question.

## 5 CONCLUSIONS AND FUTURE WORK

### 5.1 CONCLUSIONS

In this paper, we majorly studied the multi-hop reasoning tasks with the **Agent-as-tool** framework. We found that the **Agent-as-tool** could achieve a substantial improvement in the performance of the multi-hop reasoning tasks, especially in the Bamboogle dataset (Press et al., 2023). We also found that the **Agent-as-tool** could reason with less fuzzy and more structured observations, comparing with the **Search-R1 + Web Search** which would need to reason with the unprocessed web search results with fuzzy and unstructured symbols or other unrelated details.

### 5.2 LIMITATIONS AND FUTURE WORK

This paper only assigns the search tool to the agent (or in another word, the search agent) so the scope is limited to the open-domain multi-hop search tasks. While because only 1 model was provided, the dynamic assignment of the tool to the agent is not considered. Therefore in our future work more tools would be considered to be assigned to the agent, while we would also explore the dynamic assignment of the tool to the agent, in another word, make the **Planner** as a **Tool Orchestrator**.



---

## REFERENCES

- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- A. Goldie, A. Mirhoseini, H. Zhou, I. Cai, and C. D. Manning. Synthetic data generation and multi-step rl for reasoning and tool use, 2025. URL <https://arxiv.org/abs/2504.04736>.
- D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- X. Ho, A.-K. D. Nguyen, S. Sugawara, and A. Aizawa. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. *arXiv preprint arXiv:2011.01060*, 2020.
- S. Hong, X. Zheng, J. Chen, Y. Cheng, J. Wang, C. Zhang, Z. Wang, S. K. S. Yau, Z. Lin, L. Zhou, et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 3(4):6, 2023.
- A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- B. Jin, H. Zeng, Z. Yue, J. Yoon, S. Arik, D. Wang, H. Zamani, and J. Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning, 2025. URL <https://arxiv.org/abs/2503.09516>.
- G. Li, H. Hammoud, H. Itani, D. Khizbullin, and B. Ghanem. Camel: Communicative agents for” mind” exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008, 2023a.
- M. Li, Y. Zhao, B. Yu, F. Song, H. Li, H. Yu, Z. Li, F. Huang, and Y. Li. Api-bank: A comprehensive benchmark for tool-augmented llms, 2023b. URL <https://arxiv.org/abs/2304.08244>.
- X. Li, G. Dong, J. Jin, Y. Zhang, Y. Zhou, Y. Zhu, P. Zhang, and Z. Dou. Search-o1: Agentic search-enhanced large reasoning models, 2025a. URL <https://arxiv.org/abs/2501.05366>.
- X. Li, H. Zou, and P. Liu. Torl: Scaling tool-integrated rl. *arXiv preprint arXiv:2503.23383*, 2025b.
- X. Liang, J. Xiang, Z. Yu, J. Zhang, S. Hong, S. Fan, and X. Tang. Open-manus: An open-source framework for building general ai agents, 2025. URL <https://doi.org/10.5281/zenodo.15186407>.
- B. Liu, X. Li, J. Zhang, J. Wang, T. He, S. Hong, H. Liu, S. Zhang, K. Song, K. Zhu, et al. Advances and challenges in foundation agents: From brain-inspired intelligence to evolutionary, collaborative, and safe systems. *arXiv preprint arXiv:2504.01990*, 2025.
- O. Press, M. Zhang, S. Min, L. Schmidt, N. A. Smith, and M. Lewis. Measuring and narrowing the compositionality gap in language models, 2023. URL <https://arxiv.org/abs/2210.03350>.
- C. Qian, E. C. Acikgoz, Q. He, H. Wang, X. Chen, D. Hakkani-Tür, G. Tur, and H. Ji. Toolrl: Reward is all tool learning needs. *arXiv preprint arXiv:2504.13958*, 2025.
- Qwen, :, A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, H. Lin, J. Yang, J. Tu, J. Zhang, J. Yang, J. Yang, J. Zhou, J. Lin, K. Dang, K. Lu, K. Bao, K. Yang, L. Yu, M. Li, M. Xue, P. Zhang, Q. Zhu, R. Men, R. Lin, T. Li, T. Tang, T. Xia, X. Ren, X. Ren, Y. Fan, Y. Su, Y. Zhang, Y. Wan, Y. Liu, Z. Cui, Z. Zhang, and Z. Qiu. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.

- 
- Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- H. Song, J. Jiang, Y. Min, J. Chen, Z. Chen, W. X. Zhao, L. Fang, and J.-R. Wen. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. *arXiv preprint arXiv:2503.05592*, 2025.
- H. Trivedi, N. Balasubramanian, T. Khot, and A. Sabharwal. Musique: Multihop questions via single-hop question composition, 2022. URL <https://arxiv.org/abs/2108.00573>.
- L. von Werra, Y. Belkada, L. Tunstall, E. Beeching, T. Thrush, N. Lambert, S. Huang, K. Rasul, and Q. Gallouédec. Trl: Transformer reinforcement learning. <https://github.com/huggingface/trl>, 2020.
- Z. Yang, P. Qi, S. Zhang, Y. Bengio, W. W. Cohen, R. Salakhutdinov, and C. D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.
- S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao. React: Synergizing reasoning and acting in language models, 2023. URL <https://arxiv.org/abs/2210.03629>.
- J. Zhang, J. Xiang, Z. Yu, F. Teng, X. Chen, J. Chen, M. Zhuge, X. Cheng, S. Hong, J. Wang, et al. Aflow: Automating agentic workflow generation. *arXiv preprint arXiv:2410.10762*, 2024.
- Y. Zheng, S. Sun, L. Qiu, D. Ru, C. Jiayang, X. Li, J. Lin, B. Wang, Y. Luo, R. Pan, et al. Openresearcher: Unleashing ai for accelerated scientific research. *arXiv preprint arXiv:2408.06941*, 2024.
- Y. Zheng, D. Fu, X. Hu, X. Cai, L. Ye, P. Lu, and P. Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv preprint arXiv:2504.03160*, 2025.

## A THE QUALITATIVE RESULTS INSPECTION

### A.1 THE QUALITATIVE RESULTS INSPECTION OF THE **Agent-as-tool-Instruct + Web Search** AND THE **Search-R1 + Web Search**

Question: Where was the mother of Ernst I, Prince Of Hohenlohe-Langenburg born? Expected Answer: Kliczków	
Result from Search-R1	Result from Agent-as-a-tool (Ours)
<pre> &lt;think&gt; I need to find out where the mother of Ernst I, Prince of Hohenlohe-Langenburg, was born. I'll search for it. &lt;/think&gt; &lt;search&gt; mother of Ernst I, Prince of Hohenlohe-Langenburg &lt;/search&gt; &lt;information&gt; [{'result_id': 1, 'title': 'Ernst I, Prince of Hohenlohe-Langenburg', 'link': 'https://en.wikipedia.org/wiki/Ernst_I,_Prince_of_Hohenlohe- Langenburg', 'snippet': 'He married Princess Feodora of Leiningen, the only daughter of Emich Carl, 2nd Prince of Leiningen, and Princess Victoria of Saxe-Coburg-Saalfeld on 18\xa0}.... (save space for 9 more results)] &lt;/information&gt; </pre>	<pre> &lt;think&gt; I need to find out who the mother of Ernst I, Prince of Hohenlohe-Langenburg is first. &lt;/think&gt; &lt;tool_calling&gt; Who is the mother of Ernst I, Prince of Hohenlohe-Langenburg? &lt;/tool_calling&gt; &lt;obs&gt; The mother of Ernst I, Prince of Hohenlohe-Langenburg, is Countess Amalie Henriette of Solms-Baruth. For more information, you can visit the [Wikipedia page on Ernst I, Prince of Hohenlohe- Langenburg](https://en.wikipedia.org/wiki/Ernst_I,_Prince_of_H ohenlohe-Langenburg). &lt;/obs&gt; </pre>
<pre> &lt;think&gt; I found out that the mother of Ernst I, Prince of Hohenlohe- Langenburg, was Princess Feodora of Leiningen. Now I need to find out where she was born. &lt;/think&gt; &lt;search&gt; where was Princess Feodora of Leiningen born &lt;/search&gt; &lt;information&gt; [{'result_id': 1, 'title': 'Princess Feodora of Leiningen', 'link': 'https://en.wikipedia.org/wiki/Princess_Feodora_of_Leiningen', 'snippet': 'Life. Feodora was born in Amorbach, Bavaria, on 7 December 1807 to Princess Victoria of Saxe-Coburg-Saalfeld and her first husband, Emich Carl, Prince of\xa0}....', {'result_id': 2, 'title': 'Princess Feodora of Leiningen: Queen Victoria's Half Sister', 'link': 'https://owlcation.com/humanities/Feodora-of- Leiningen-Queen-Victorias-Half-Sister', 'snippet': 'Feodora was born Princess Anna Feodora Auguste Charlotte Wilhemina of Leiningen on the 7th of December 1807 in Amorbach, then in Leiningen and\xa0}....'.... (save the space for 8 more search results)] &lt;/information&gt; </pre>	<pre> &lt;think&gt; Now that I know the mother's name, I need to find out where she was born. &lt;/think&gt; &lt;tool_calling&gt; Where was Countess Amalie Henriette of Solms-Baruth born? &lt;/tool_calling&gt; &lt;obs&gt; Countess Amalie Henriette of Solms-Baruth was born in Kliczków on January 30, 1768. &lt;/obs&gt; &lt;think&gt; The information provided states that Countess Amalie Henriette of Solms-Baruth, the mother of Ernst I, Prince of Hohenlohe- Langenburg, was born in Kliczków. &lt;/think&gt; &lt;answer&gt; <b>Kliczków (Correct Answer)</b> &lt;/answer&gt; </pre>
<pre> &lt;think&gt; I found out that Princess Feodora of Leiningen was born in Amorbach, Bavaria, Germany. Now I can provide the answer. &lt;/think&gt; &lt;answer&gt; <b>Amorbach (Wrong Answer)</b> &lt;/answer&gt; </pre>	

Figure 3: The **Agent-as-tool-Instruct** could reason with less fuzzy and more structured observations, comparing with the **Search-R1 + Web Search** which would need to reason with the unprocessed web search results with other unrelated details. Search-R1 was misled by the unprocessed web search results to reason with the wrong answer for the second hop (Princess Feodora of Leiningen), while the **Agent-as-tool-Instruct** has applied the agent to preprocess the web search results and return the correct answer (Countess Amalie Henriette of Solms-Baruth) for the second hop.

## A.2 THE QUALITATIVE RESULTS INSPECTION OF THE **Agent-as-tool-Instruct + Web Search** AND THE **Agent-as-tool-Base + Web Search**

Question: Where did Edward Dunn (Bishop)'s father graduate from? Expected Answer: Corpus Christi College, Cambridge	
Result from Agent-as-tool-Base	Result from Agent-as-tool-Instruct
<p>&lt;think&gt; I need to find Edward Dunn's father's educational background. &lt;/think&gt;</p> <p>&lt;tool_calling&gt; Can you provide information on Edward Dunn's father's education? &lt;/tool_calling&gt;</p> <p>Planner try to ask 2 hops simultaneously</p> <p>&lt;obs&gt; I wasn't able to find specific information about Edward Dunn's father's education. However, I did find information related to Edward Dunn himself, including his education and achievements. If you need information on a different topic or a specific Edward Dunn, please let me know! &lt;/obs&gt;</p> <p>The agent was unable to answer</p> <p>&lt;think&gt; The observation did not provide the specific information I was looking for. I need to find a more direct way to get the information about Edward Dunn's father's education. &lt;/think&gt;</p> <p>&lt;tool_calling&gt; Can you provide information on the educational background of Edward Dunn's father? &lt;/tool_calling&gt;</p> <p>Planner asked 2 hops in the meantime in another manner</p> <p>&lt;obs&gt; I found some relevant information regarding Edward Dunn's family background, but specific details about his father's educational background are limited. Here are some notable findings: 1. <b>Edward T. Dunn</b> - He is noted to have studied theology and earned a master's degree in history and a PhD, but this relates to him rather than his father. [Source](https://blogs.canisius.edu/the-dome/2013/03/26/canisius-mourns-the-loss-of-rev-edward-t-dunn-s-j/). 2. <b>Bishop Edward Dunn</b> - His father, Andrew Hunter Dunn, was an education graduate from <u>Marlborough and Pembroke College, Cambridge</u>. This lineage could provide context about the educational background in the Dunn family. [Source](https://en.wikipedia.org/wiki/Edward_Dunn_(bishop)). 3. Other sources mainly discuss different individuals named Edward Dunn or their obituaries, without providing specific information about his father's education. If you would like more specific information or about a different Edward Dunn, please clarify! &lt;/obs&gt;</p> <p>The agent cannot get exact information</p> <p>&lt;think&gt; The observation provided the specific information I needed. Edward Dunn's father, Andrew Hunter Dunn, graduated from Marlborough and Pembroke College, Cambridge. &lt;/think&gt;</p> <p>&lt;answer&gt; <b>Marlborough and Pembroke College, Cambridge (Wrong Answer)</b> &lt;/answer&gt;</p>	<p>&lt;think&gt; I need to find information about Edward Dunn (Bishop)'s father. &lt;/think&gt;</p> <p>&lt;tool_calling&gt; Can you provide information about Edward Dunn (Bishop)'s father? &lt;/tool_calling&gt;</p> <p>Learned to Ask the prerequisite question before the next hop</p> <p>&lt;obs&gt; Edward Dunn, the bishop, was born into an ecclesiastical family. His father was <u>Andrew Hunter Dunn</u>, who at one time served as the Bishop of Quebec. For more details, you can read the full biography on [Wikipedia](https://en.wikipedia.org/wiki/Edward_Dunn_(bishop)). &lt;/obs&gt;</p> <p>The agent could answer this question easily</p> <p>&lt;think&gt; From the observation, I can see that Edward Dunn's father was Andrew Hunter Dunn, who was the Bishop of Quebec. However, the observation does not specify where Andrew Hunter Dunn graduated from. I need to find this information. &lt;/think&gt;</p> <p>&lt;tool_calling&gt; Can you provide information about where Andrew Hunter Dunn graduated from? &lt;/tool_calling&gt;</p> <p>Asked the next-hop question</p> <p>&lt;obs&gt; Andrew Hunter Dunn graduated from <b>Corpus Christi College, Cambridge</b>, where he earned his BA as the 29th Wrangler in 1863. For more detailed information, you can visit his [Wikipedia page](https://en.wikipedia.org/wiki/Andrew_Hunter_Dunn). &lt;/obs&gt;</p> <p>The agent was able to answer the question easily</p> <p>&lt;think&gt; Now I have the information that Andrew Hunter Dunn, the father of Edward Dunn (Bishop), graduated from Corpus Christi College, Cambridge. &lt;/think&gt;</p> <p>&lt;answer&gt; <b>Corpus Christi College, Cambridge (Correct Answer)</b> &lt;/answer&gt;</p>

Figure 4: The **Agent-as-tool-Instruct + Web Search** could correctly decompose the question to identify the first hop and the second hop, comparing with the **Agent-as-tool-Base + Web Search** which barely decompose the question and try to ask about the whole question in another manner, i.e. the **Agent-as-tool-Base + Web Search** was not able to reason the whole multi-hop question correctly, while the **Agent-as-tool-Instruct + Web Search** could decompose the question to identify the first hop of the question to be answered by the agent then proceed to the second hop.