

Scaling Long-Horizon LLM Agent via Context-Folding

Weiwei Sun^{1,2,*}, Miao Lu^{1,3,*}, Zhan Ling¹, Kang Liu¹, Xuesong Yao¹, Yiming Yang²,
Jiecao Chen^{1,†}

¹ByteDance Seed, ²Carnegie Mellon University, ³Stanford University

*Work done at ByteDance Seed, †Corresponding authors

Abstract

Large language model (LLM) agents are fundamentally constrained by context length on long-horizon tasks. We introduce Context-Folding, a framework that empowers agents to actively manage their working context. An agent can procedurally branch into a sub-trajectory to handle a subtask and then fold it upon completion, collapsing the intermediate steps while retaining a concise summary of the outcome. To make this behavior learnable, we develop an end-to-end reinforcement learning framework FoldGRPO with specific process rewards to encourage effective task decomposition and context management. On complex long-horizon tasks (Deep Research and SWE), our folding agent matches or outperforms the ReAct baselines while using an active context 10× smaller and significantly outperforms models that rely on summarization-based context management.

Date: October 15, 2025

Correspondence: Weiwei Sun at sunnweiwei@gmail.com, Jiecao Chen at jiecao.chen@bytedance.com

Project Page: <https://context-folding.github.io/>

1 Introduction

Large language model (LLM) agents have shown remarkable capabilities in tackling complex, long-horizon problems that require extensive interaction with an environment, such as deep research [8, 12, 17, 21, 32] and agentic coding [3, 11, 31]. The length of tasks agents can complete is argued to be *growing exponentially, with a doubling time of about 7 months* [20].

However, scaling LLM agents to even longer horizons is fundamentally constrained by the design of agentic frameworks [35]. These frameworks linearly accumulate the entire interaction history (reasoning, tool calls, and observations) into a single, ever-expanding context, which incurs long-context challenges as horizons scale: (1) degraded performance, as LLMs struggle to utilize relevant information in exceedingly long contexts [14, 18, 28]; and (2) poor efficiency, stemming from the quadratic scaling of attention mechanisms and the growing overhead of managing the KV-cache [13].

Existing approaches to scale long-horizon LLM agents largely fall into two classes: (1) *Summary-based methods*, which trigger a post-hoc summarization stage when the working context is full [1, 19, 24, 34, 38, 43]. While this compresses the context, it can abruptly disrupt the agent’s working context and reasoning flow, which may lead to sub-optimal results. (2) *Multi-agent systems*, which distribute tasks across specialized agents to manage context length [2, 33, 40, 41]. Yet, these systems typically depend on handcrafted, problem-specific workflows that are difficult to generalize and resist end-to-end optimization.

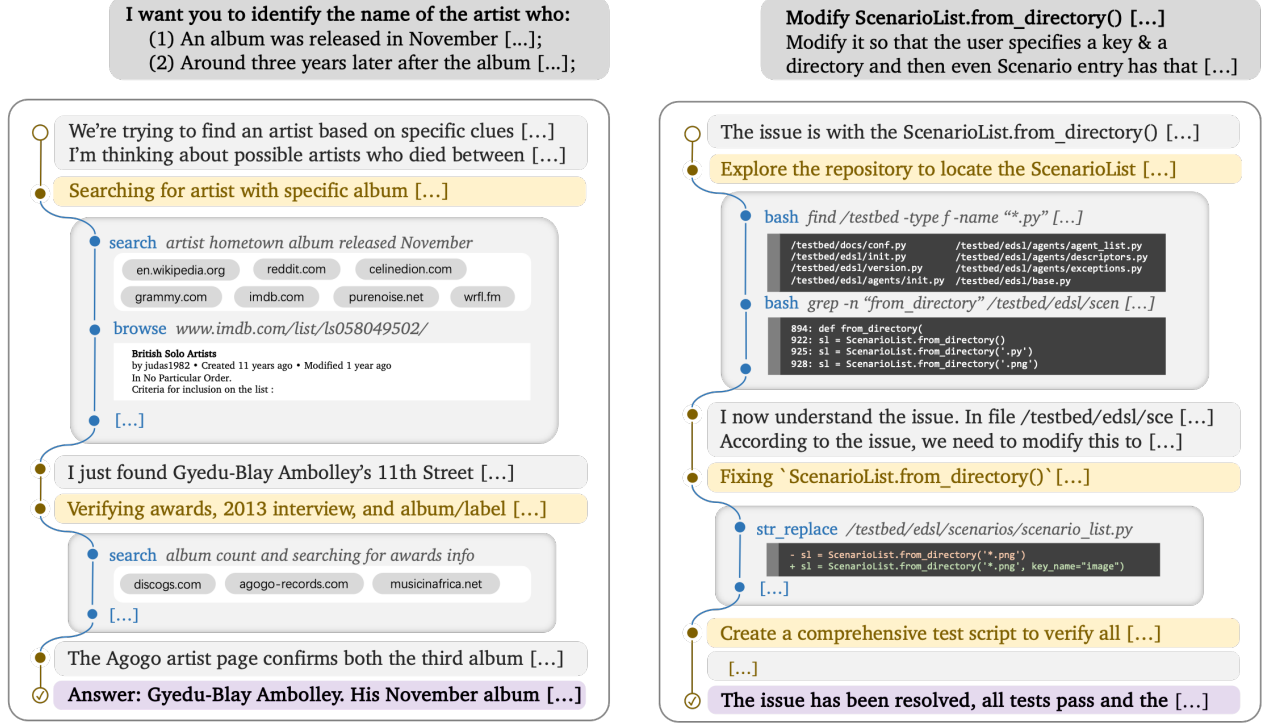


Figure 1 Examples of context folding in long-horizon tasks: deep research (left) and agentic coding (right).

In this paper, we propose **Context Folding**: an agentic mechanism that allows the model to actively manage its working context. Specifically, the agent manages its context using two special actions: (i) a **branch** action to create a temporary sub-trajectory for a localized subtask; and (ii) a **return** action to summarize the outcome and rejoin the main thread, after which the intermediate steps within the branch are “folded”—removed from the context—leaving only a concise summary from the **return** call. Figure 1 illustrates this process on deep research and agentic coding tasks, where the agent offloads token-intensive actions (e.g., web search or codebase exploration) into branches and preserves only key findings and insights for high-level reasoning. Compared with existing methods, context folding enables an agentic approach to active context management, where the agent’s short-term context remains undisrupted and long-term context is automatically managed.

Based on the context-folding framework, we propose a novel end-to-end reinforcement learning algorithm for training LLM agents on complex, long-horizon tasks. The key innovation is **FoldGRPO**, which augments the standard GRPO by incorporating (i) dynamic folded LLM contexts and (ii) dense, token-level process rewards that directly guide context folding behavior. Specifically, our RL algorithm teaches the model how to effectively decompose a problem into localized sub-tasks for branching, guided by an *Unfolded Token Penalty* that discourages token-heavy operations in the main context. Furthermore, it learns to maintain focus within a sub-task via an *Out-of-Scope Penalty*, and to preserve crucial information in its summaries to aid the final objective. By mastering these skills, the agent can handle vastly longer interaction histories, allowing our framework to scale the agent’s effective horizon and improve overall system efficiency.

We evaluate our approach on two long-horizon benchmarks, BrowseComp-Plus [6] and SWE-Bench Verified [11], where our agent achieves strong performance with remarkable efficiency. Despite using a compact 32K active token budget managed with maximum of 10 branches, our agent, the **Folding Agent**, achieves pass@1 scores of 62.0% and 58.0% respectively, surpassing baselines that require a massive 327K context window and significantly outperforming methods based on context summarization. The effectiveness of our method is rooted in reinforcement learning, which provides absolute improvements of 20.0% on BrowseComp-Plus and 8.8% on SWE-Bench. Further analysis reveals that our agent learns to invoke more tool calls and generate longer outputs to handle complex problems, and can scale to larger token budgets at inference time to tackle

even more challenging tasks. Together, these results indicate that learning to *actively manage* context, rather than merely extending or heuristically compressing it, is a principled path toward scalable long-horizon agency.

In summary, our contributions are threefold:

- (i) We introduce **Context Folding**, a mechanism that enables agents to actively manage their context and mitigate the challenge of linear history growth.
- (ii) We present **FoldGRPO**, a reinforcement learning framework with dynamic folded LLM contexts and dense process rewards that trains agents to effectively acquire the capability of context folding.
- (iii) We demonstrate promising performance on long-horizon benchmarks, highlighting our approach as a scalable and extensible path toward stronger LLM agents.

2 Methodology

2.1 Vanilla Formulation

Given a question q , an agent generates a multi-turn interaction trajectory denoted as

$$\tau := (a_1, o_1, a_2, o_2, \dots, a_T, o_T),$$

where a_i is the LLM output at step i (including *reasoning* and *tool call*), and o_i is the corresponding tool-call result. The vanilla ReAct-style agent [35] models the interaction as following,

$$p_{\theta}^{\text{ReAct}}(\tau | q) = \prod_{i \in [T]} \pi_{\theta}(a_i | q, (a_1, o_1, \dots, a_{i-1}, o_{i-1})),$$

which appends the entire interaction history to the context at each time of LLM generation. However, in long-horizon agentic tasks like deep research and agentic coding, τ can accumulate rapidly due to extensive interactions and become prohibitively long which exceeds the working context limit. Also, when the context is expanding, the reasoning and instruction following capability of the model may drop, posing further challenges for the agent to complete the long-horizon task.

2.2 Our Method: Context Folding

To address the challenge, we introduce *context folding*, a mechanism that allows the agent to *actively manage its working context via branching and folding*. Specifically, we design two tools that the agent can call for context management. Starting from a main thread to solve question q , it can:

- (i) **branch**(*description*, *prompt*): *branch from main thread to use a separate working context to complete a sub-task q' for solving q . Here *description* is a brief summary of the sub-task, and *prompt* is a detailed instruction for this branch. The tool returns a template message indicating that the branch was created.*
- (ii) **return**(*message*): *fold the context generated in this branch and return to the main thread. The *message* describes the outcome of this branch. Upon calling this tool, the agent context then switches back to the main thread, appended with the templated *message* from the branch.*

With these two tools, the agent can actively manage its context by (i) branching a separate working context to solve an independent sub-task, and (ii) folding the intermediate steps in the branch, and resuming back to the main thread by appending only the result of the branch. To put it formal, the context-folding agent is modeled as following,

$$p_{\theta}^{\text{Context Fold}}(\tau | q) := \prod_{i \in [T]} \pi_{\theta}(a_i | q, \mathcal{F}(\tau_{<i})). \quad (1)$$

Here $\tau_{<i} = (a_1, o_1, \dots, a_{i-1}, o_{i-1})$ denotes the complete history of all the action-observation pairs before step i , \mathcal{F} is the context manager that folds the interaction history between **branch** and **return** tool calls. We

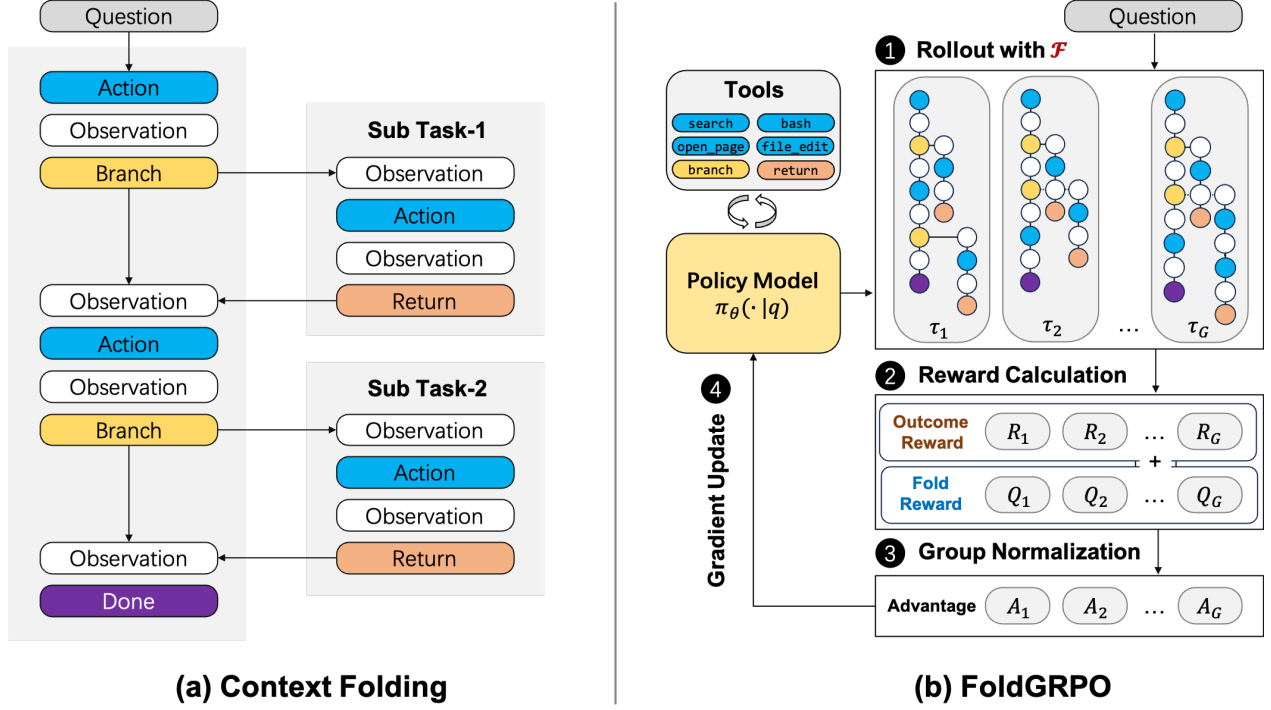


Figure 2 (a) Context Folding: a mechanism that enables the agent to actively manage its context through branching and return. **(b) FoldGRPO:** end-to-end optimization of context folding agent.

illustrate the process using the following example, where the context manager folds all the action-observation pairs in previous branches:

$$\begin{aligned}
 &\mathcal{F}(a_1, o_1, \underbrace{a_2, o_2, a_3, o_3, a_4, o_4}_{\text{branch 1}}, \underbrace{a_5, o_5, a_6, o_6, a_7, o_7, a_8, o_8}_{\text{branch 2}}, a_9, o_9, a_{10}, o_{10}) \\
 &\quad \rightarrow (a_1, o_1, a_2, o_4, a_5, o_8, a_9, o_9, a_{10}, o_{10}),
 \end{aligned}$$

so the segments between a_2 and a_4 and between a_5 and a_8 are folded.

Inference efficiency. During inference, the agent manages a context KV-cache: when `return` action is called, it rolls back the KV-cache to the corresponding `branch` position, where the context prefix matches that before calling the `branch` action. This makes our context folding approach efficient in terms of inference.

Instantiation: plan-execution. To instantiate context folding for long-horizon tasks, we adopt a *plan-execution* framework, where the agent alternates between two states: (i) *Planning State*: The agent performs high-level reasoning in the main thread, decomposes the task, and decides when to initiate a branch for a sub-task. In this state, token-intensive tool use is discouraged to keep the main context focused on high-level strategies. (ii) *Execution State*: The agent operates within an active branch to complete its assigned sub-task. To maintain a clear structure and prevent nested complexity, creating new branches is disabled while in execution state.

2.3 FoldGRPO: End-to-End RL for Context-Folding Agent

To optimize the context folding agent, in this section, we introduce an end-to-end RL training framework, namely, Folded-context Group Relative Policy Optimization (FoldGRPO). FoldGRPO jointly optimizes the entire interaction trajectory including the main thread and those sub-task branches, while it folds the rollout history according to the context folding modeling (1) to maintain a compact working context during training. Moreover, FoldGRPO features a novel *process reward design* to efficiently guide the training of the branching behavior of the agent. We first introduce the overall algorithm design in Section 2.3.1 and we present the process reward design in Section 2.3.2.

2.3.1 Overall Algorithm Design

In each training step of FoldGRPO, for task q from training dataset \mathcal{D} , G trajectories $(\tau_1, \tau_2, \dots, \tau_G)$ are sampled from the old policy π_{old} according to the context folding model (1). Each complete trajectory, e.g., $\tau_i = (a_{i,1}, o_{i,1}, \dots, a_{i,T}, o_{i,T})$, is a sequence of tokens defined as $\tau_i = [\tau_{i,1}, \dots, \tau_{i,|\tau_i|}]$. Each trajectory τ_i has a final reward $R_i \in \{0, 1\}$, following the recipe of RL from verifiable rewards (RLVR).

Learning objective. The learning objective of FoldGRPO is defined as:

$$\mathcal{J}_{\text{FoldGRPO}} = \mathbb{E}_{\substack{q \sim \mathcal{D}, \\ \{\tau_i\}_{i=1}^G \sim \pi_{\text{old}}(\cdot|q)}} \left[\frac{1}{\sum_{i=1}^G |\tau_i|} \sum_{i=1}^G \sum_{t=1}^{|\tau_i|} \min \left\{ r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip}(r_{i,t}(\theta), 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}}) \hat{A}_{i,t} \right\} \right],$$

where the importance sampling ratio and the group relative advantage estimator [25] are given by

$$r_{i,t}(\theta) = \frac{\pi_{\theta}(\tau_{i,t} | q, \mathcal{F}(\tau_{i,<t}))}{\pi_{\theta_{\text{old}}}(\tau_{i,t} | q, \mathcal{F}(\tau_{i,<t}))} \cdot \mathbf{1}_{\tau_{i,t}}^{\text{LLM}}, \quad \hat{A}_{i,t} = \frac{\text{clip}(R_i + Q_{i,t}, 0, 1) - \text{mean}(\{R_i\}_{i=1}^G)}{\text{std}(\{R_i\}_{i=1}^G)}.$$

Here $\mathbf{1}_{\tau_{i,t}}^{\text{LLM}}$ ensures that only those LLM generated tokens are optimized and the tokens from tool observations are masked. In the following, we explain two key features of FoldGRPO highlighted in red.

- (i) **Context folding.** Unlike vanilla multi-turn LLM RL algorithms that append the entire interaction history to context when optimizing the policy, FoldGRPO applies context manager $\mathcal{F}(\cdot)$ to the history $\tau_{i,<t}$ which folds the context for token $\tau_{i,t}$ based on the branch-return actions
- (ii) **Process reward signal.** In the calculation of advantage $\hat{A}_{i,t}$, a token-level process reward $Q_{i,t}$ is added to regularize the model’s branch-return behavior, which is detailed in the next section.

2.3.2 Process Reward Design

In RLVR, the agent is typically optimized through a standard binary *outcome reward* based on task success or failure. However, we empirically observe that this sparse reward signal is insufficient for learning effective context folding. Specifically, two critical failure modes emerge: (i) The agent fails to plan strategically, leaving token-intensive operations unfolded in the main context, which quickly exhausts the available token budget. (ii) The agent struggles with proper branch management, often failing to return from a sub-branch after a sub-task is completed and instead continuing the subsequent work within that same branch. To effectively optimize the folding agent, we introduce token-level process rewards separately to main-trajectory tokens and branch-trajectory tokens.

Unfolded token penalty. When total context length of the main thread exceeds 50% of the working context limit, we apply $Q_{i,t} = -1$ to all the tokens in the main thread, except those tokens in the turns that create a branch. This penalizes the agent for performing token-heavy actions outside a branch in the main thread, and encourages the agent to perform those actions in separate branches.

Out-scope penalty. For each branch, we employ GPT-5-nano to judge — based on the branch prompt and the returned message — whether the agent has conducted actions outside the specified sub-tasks. If so, we apply $Q_{i,t} = -0.2$ to all the tokens in this branch to penalize such out of scope behavior. This encourages the agent to only perform the exact sub-task given to the current branch.

Failure penalty. We apply $Q_{i,t} = -1$ to all the tokens in a failed tool call turn. In all other cases, $Q_{i,t} = 0$.

2.4 How does Context Folding Connect to Other Methods?

Relationship to multi-agent systems. Context folding can be interpreted as a specific formulation of a general multi-agent system, where the main agent delegates sub-tasks to sub-agents. Compared to popular multi-agent systems [9], our design differs in the following ways: (i) Context folding does not adopt predefined sub-agents; instead, sub-agents are created by the main agent on the fly; (ii) All the agents share the same context prefix, making it KV-cache friendly, (iii) The main and the sub agents interleave rather than operating in parallel.

Relationship to context-summarization-based method. Compared with heuristic summarization-based context management [21, 38], which discards details at arbitrary points, context folding can be viewed as a learnable summarization mechanism aligned with sub-task boundaries. This ensures that reasoning is preserved during execution and is only compacted once its utility is realized.

3 Experiment

3.1 Datasets

We conduct experiment on two representative long-horizon agent tasks: deep research, and agentic software engineering:

Deep Research. We use BrowseComp-Plus (BC-Plus) [6], which supplements the original BrowseComp data with a verified corpus. We use Qwen3-Embed-8B as the retriever. Since the quality of training data is crucial for the BrowseComp task but existing datasets are typically not open-sourced [15, 24], we split BrowseComp-Plus into training and evaluation sets to decouple the effect of data distribution. Our split includes 680 instances for training and 150 for evaluation. For deep research, the tools are `search(query, topk)` and `open_page(url)`, and the reward is based on official LLM-based judge [6].

Agentic SWE. We use SWE-Bench Verified (SWEB-V) [11] as the evaluation set. To collect training data, we roll out the baseline agent¹ eight times on a subset of the open-source datasets SWE-Gym [23] and SWE-Rebench [4], and retain the instances where the success rate is between 0 and 87.5%, resulting in 740 instances. In SWE, the tools are `execute_bash`, `str_replace_editor`, and `think` [31], and the reward is based on running unit test in instance-specific sandbox environment.

We classify test instances for both tasks into three difficulty levels: *easy*, *medium*, and *hard*. For BrowseComp-Plus, classification is determined by running a ReAct agent 8 times on each instance. An instance is labeled *easy* if its `acc@8` score is $\geq 87.5\%$, *hard* if its score is 0% , and *medium* otherwise, resulting in 50 instances for each level. For SWE-Bench Verified, classification is based on the original dataset’s time-to-resolve metric: *easy* (≤ 15 min, 194 instances), *medium* (15 min – 1 hour, 261 instances), and *hard* (≥ 1 hour, 45 instances).

See Appendix B for the details of system prompt of each datasets.

3.2 Implementation

We use Seed-OSS-36B-Instruct² as the base LLM and conduct RL training on it. For RL training, we build on VeRL and set the rollout batch size to 32, group size to 8, ppo batch size of 128, learning rate to 1×10^{-6} , no KL term, clip high to 0.28, and clip low to 0.2. We employ asynchronous rollout with a maximum off-policy step of 5. During training, we implement the context folding operation \mathcal{F} by constructing separate causally conditioned contexts for each branch to improve training efficiency (See Appendix A for more details.). We train model for 50 steps (about 2 epochs). For the fold agent, we set the LLM maximum context length to 32,768. We allow up to 10 branches, resulting in a theoretical maximum of 327,680 tokens. During inference we employ greedy decoding (i.e, temperature = 0).

3.3 Baselines

We compare against the following baselines:

ReAct Agent [36], which keeps all context. We consider different context lengths for comparison: (a) *short-context*, which has 32,768 tokens, equivalent to our context length; (b) *medium-context*, which has intermediate lengths, e.g., 65,536 and 131,072; (c) *long-context*, which has 327,680 tokens, equivalent to our maximum total token cost.

Summary Agent [34, 38], which invokes a summary when the context is full. We set the maximum context length to 32,768 and allow for 10 summary session for a fair comparison.

¹Seed-OSS-36B-Instruct with OpenHands and a response length of 65,536.

²<https://huggingface.co/ByteDance-Seed/Seed-OSS-36B-Instruct>

Model	Peak Length	Max #Token	BrowseComp-Plus		SWE-Bench Verified	
			Pass@1	Tool Calls	Pass@1	Tool Calls
ReAct Agent with 100B+ LLM						
GPT-5	327K	327K	0.793	14.2	0.718	42.6
GPT-4.1	327K	327K	0.640	5.6	0.486	28.7
DeepSeek-V3.1	327K	327K	0.613	10.6	0.610	53.2
GLM-4.5-Air	327K	327K	0.566	11.1	0.576	51.2
Qwen3-235B-A22B	327K	327K	0.560	12.8	0.344	32.1
ReAct Agent						
Seed-OSS-36B	32K	32K	0.286 (-19.2)	3.8	0.436 (-11.6)	25.8
+ RL (GRPO)	32K	32K	0.446 (-3.2)	5.5	0.480 (-7.2)	27.8
Seed-OSS-36B ^ψ	327K	327K	0.478 (+0.0)	10.8	0.552 (+0.0)	49.5
+ RL (GRPO)	327K	327K	0.540 (+6.2)	10.2	0.574 (+2.2)	55.4
Summary Agent						
Seed-OSS-36B	32K	32K × 10	0.386 (-9.2)	17.4	0.488 (-6.4)	77.0
+ RL (GRPO)	32K	32K × 10	0.527 (+4.9)	18.0	0.550 (-0.2)	74.9
Folding Agent (Ours)						
Seed-OSS-36B	32K	32K × 10	0.420 (-5.8)	12.9	0.492 (-6.0)	72.8
+ RL (GRPO)	32K	32K × 10	0.567 (+8.9)	16.0	0.564 (+1.2)	79.5
+ RL (FoldGRPO)	32K	32K × 10	0.620 (+14.2)	19.2	0.580 (+2.8)	96.5

Table 1 Performance on BrowseComp-Plus (N=150) and SWE-Bench Verified (N=500). Boldface indicates the best-performing 36B models. Numbers in parentheses indicate improvement or reduction compared to 327K ReAct agent Seed-OSS-36B baseline^ψ.

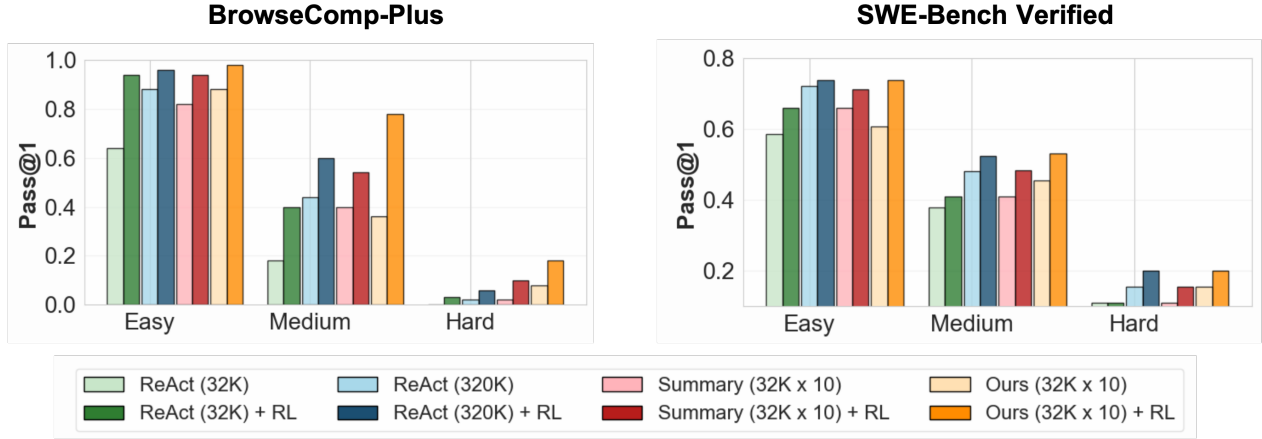


Figure 3 Agent performance on different data difficulty group. RL training yields consistent performance gains across easy, medium, and hard instances.

For both two baselines, we employ the same base model (i.e., Seed-OSS-36B-Instruct), data, infrastructure, and hyperparameters for RL training. In addition to these directly comparable baselines, we compare our method against previous closed-source and open-source systems on both tasks, including GPT-5, GPT-4.1, DeepSeek-V3.1 (2509), GLM-4.5-Air, and Qwen3-235B-A22B-Instruct-2507.

4 Experimental Results

4.1 Main Results

Table 1 summarizes our main results on the BrowseComp-Plus and SWE-Bench Verified datasets. Our findings highlight the critical role of reinforcement learning in unlocking the capabilities of context folding.

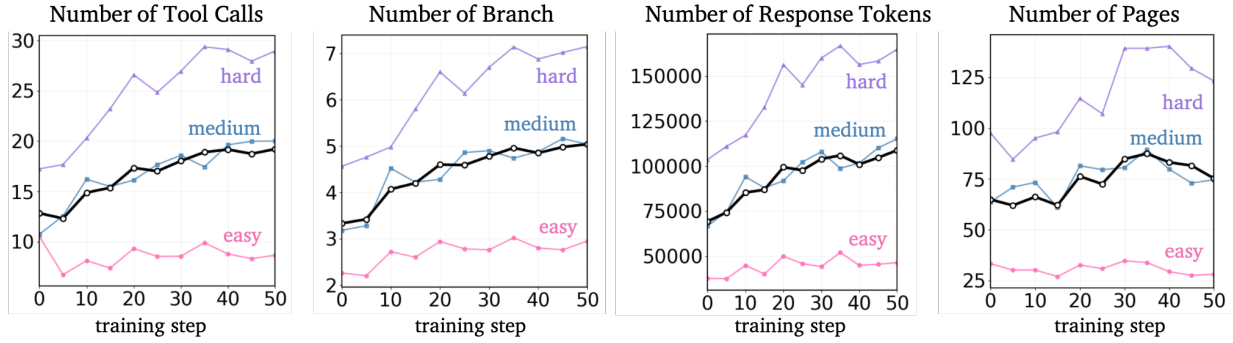


Figure 4 With RL training, we observe an increase in the number of tool calls, branching behavior, total number of tokens, and the number of searched pages.

Initially, without performing RL, the context folding agent already surpasses the 32K-context ReAct and context summarization baselines, though it does not yet match the performance of the long-context ReAct agent. After RL training, our agent’s performance improves significantly, with a *pass@1* of **0.620 on BrowseComp-Plus (+20%)** and **0.580 on SWE-Bench Verified (+8.8%)**. Our agent not only outperforms all baselines, including the long-context ReAct agent with same 327K max length, but also achieves performance comparable to agents built on much larger 100B+ parameter models.

Further analysis reveals two key insights. First, an ablation study confirms that our proposed **FoldGRPO is crucial**, yielding significantly better performance than the baseline GRPO algorithm (eg, +7.7% on BrowseComp and +1.6% on SWE-Bench). Second, the performance gains correlate with an increased frequency of tool calls, which RL training further encourages. This suggests our framework enables the agent to conduct a more thorough exploration of the environment to discover more robust solutions.

4.2 Performance by Task Difficulty

Figure 3 breaks down agent performance by task difficulty, comparing scores before and after reinforcement learning. The results clearly show that RL training yields consistent performance gains across easy, medium, and hard instances. Most notably, the improvements are significantly larger for the medium and hard subsets. This finding underscores our agent’s enhanced capability to handle complex problems that require more sophisticated long-context management.

Figure 4 shows the agent’s learning dynamics during RL training on BrowseComp-Plus. As training progresses, the agent steadily increases its tool calls, branch creation, response tokens, and number of pages searched. This growth is most pronounced on harder instances. For example, on the hard subset, response length rises from about 100K to over 160K tokens. These results show that the agent learns to allocate more interaction and computation to complex problems, adopting a more adaptive and effective problem-solving strategy.

4.3 Ablation of RL Algorithm

To understand how our proposed FoldGRPO shapes agent behavior, we analyze the key statistics in Table 2. These metrics include the task completion rate within the context limit (Finish), main trajectory length (Main Len), the accuracy of sub-trajectories staying on-topic (Scope), and the number of branches created (# Branch). We can see that, training with a standard GRPO baseline produces poor behaviors: agents show a lower Finish rate, generate overly long trajectories, and lose focus in sub-tasks, reflected in reduced Scope accuracy. This indicates a failure to manage context effectively.

By contrast, our FoldGRPO corrects these issues. It encourages focused branching, sharply boosting both Scope accuracy and Finish rate. Most notably, it cuts the main trajectory to about 8K tokens while processing over 100K in total—achieving over 90% context compression and demonstrating the agent’s ability to condense long interactions into a compact, useful history.

	BrowseComp-Plus				SWE-Bench Verified			
	Finish	Main Len	Scope	# Branch	Finish	Main Len	Scope	# Branch
Folding Agent (Seed-OSS-36B)	0.806	12,195	0.774	3.51	0.781	47,475	0.473	3.05
+ RL (GRPO)	0.738	22,285	0.762	3.88	0.612	48,908	0.419	3.80
+ RL (FoldGRPO)	0.935	7,752	0.895	4.98	0.962	8,885	0.754	5.90

Table 2 Model behavior statistics of different optimization methods. FoldGRPO encourages focused branching and condensed main context, boosting both scope accuracy and finish rate.

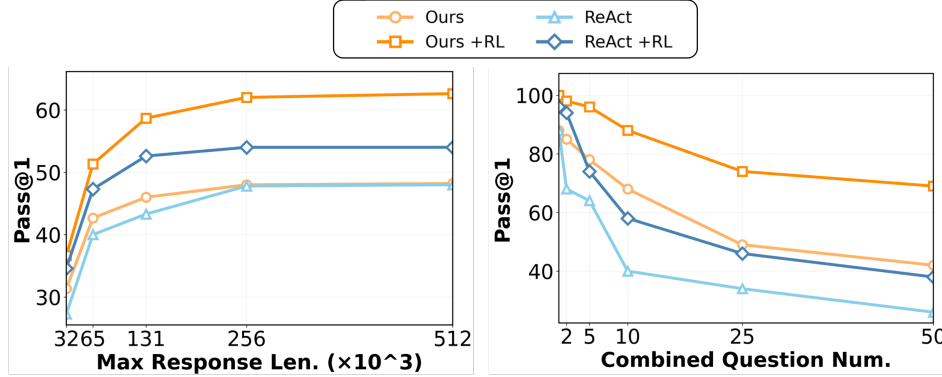


Figure 5 Left: Pass@1 vs. agent max context length. **Right:** Pass@1 vs. number of combined questions. Multiple easy questions are combined into a single harder question to increase problem complexity; a higher number of combined questions indicates more required actions and a longer context to answer them correctly. See Section 4.4.2 for details.

4.4 Performance by Context Length

4.4.1 Effect of Context Length

To examine how performance scales with context length, we evaluated our method on BrowseComp while varying the number of branches from 0 to 16. As shown in Figure 5 (left), our method consistently surpasses ReAct, and reinforcement learning provides further gains. However, performance plateaus beyond 320K tokens because most task instances are already completed, and additional context provides limited benefit.

4.4.2 Effect of Task Complexity

Following Zhou et al. [43], we increase task complexity by combining multiple questions into a single compound query that the agent must answer *in one session* (see Figure 6 for an example). Figure 5 (right) shows the results for tasks with 1 to 50 combined questions. For this setting, we allow unlimited branching and set the context limit for ReAct to 1M tokens. As task complexity increases, the benefit of context folding becomes more apparent, demonstrating strong length generalization. Notably, although our agent was trained on tasks requiring at most 10 branches, it adaptively uses an average of 32.6 branches to solve tasks with 50 questions.

4.5 Further Analysis

4.5.1 Case Study

Figure 7 shows qualitative examples of our context folding agent on BrowseComp-Plus. Given a query about finding a research publication with specific conditions, the agent first explores the high-level topic and identifies a candidate. It then searches to verify conditions, gaining key insights but failing to confirm all requirements. Next, it expands the search scope and finds the correct answer. In this process, 4 branches compress the full 107K-token context to just 6K.

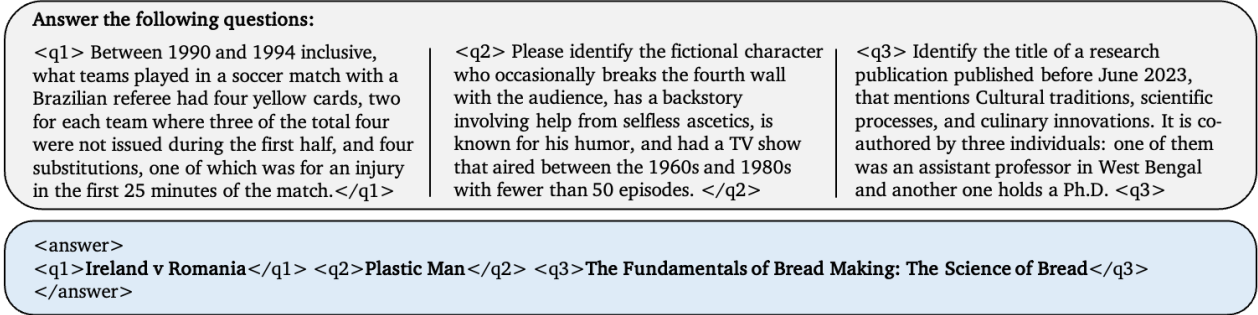


Figure 6 An example of the model’s input and output for the combined-questions experiment described in Section 4.4.2. In this example, 3 easy questions are combined to form a harder question.

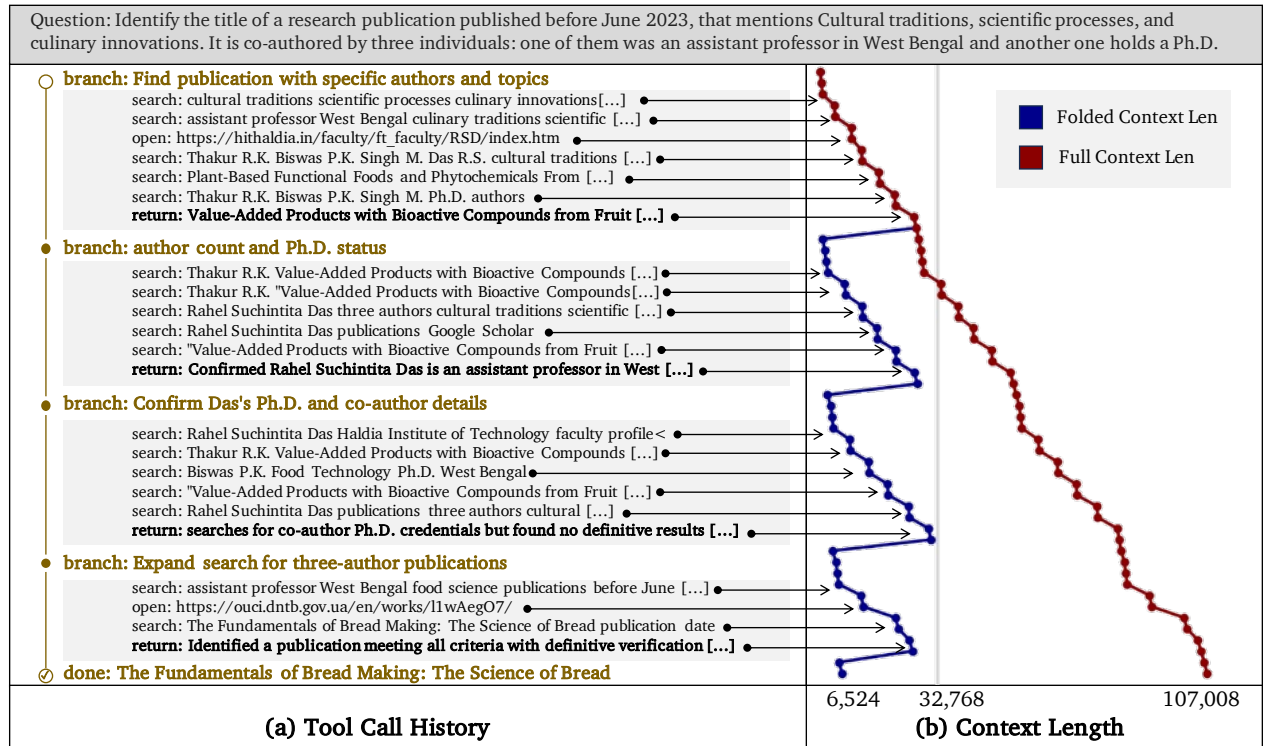


Figure 7 Example of an agent’s tool call history and context length. on BrowseComp-Plus.

4.5.2 Training Speed

Figure 8 shows the stepwise average time for rollout and for each training step. We observe that the 327K ReAct model requires a longer training time per step. Note that we employ async rollout (Appendix A.2), and the rollout time shown here measures only the main thread’s time cost on rollout.

4.5.3 Parallel Branching

Whether the folding agent can benefit from parallel branching — i.e., creating multiple sub-branches that run simultaneously — remains an open question. We experimented on BrowseComp-Plus by training an agent that utilizes parallel branching under the same setup as the single-branch agent. The parallel-branch version achieved a 0.6133 Pass@1 on BrowseComp-Plus, outperforming the baseline but performing similarly to the single-branch version. Moreover, after training, the parallel-branch agent created about 2.3 parallel branches on average and read more web pages (110 vs. 80 for the single-branch version). However, it did not achieve

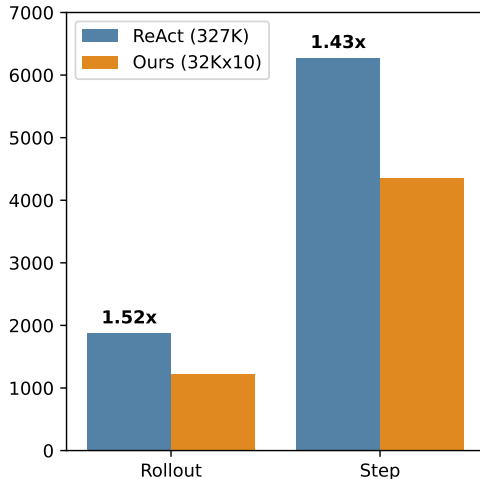


Figure 8 Training time cost. The figure shows the stepwise average time for rollout and for each training step.

a higher score, possibly because the task characteristics are more depth-first in nature. Other tasks with a breadth-first structure (eg WideSearch [33]) may be more promising for studying parallelism in LLM agents.

5 Related Work

The rapid evolution of LLM agents is driven by a push toward greater autonomy in complex, long-horizon tasks [11, 16, 20, 22, 42]. Built on agentic architectures that integrate planning, memory, and tool use [30], research has advanced from simple sequential reasoning to dynamic, multi-path strategies for exploration and problem-solving [5, 10, 26, 37]. Yet this progress has revealed a key bottleneck: the finite and costly nature of an agent’s working context [1, 35].

Context management strategies fall into two main paradigms: context summarization, where agents offload and retrieve information from external memory stores [27, 29, 34, 38, 43], and multi-agent collaboration, where tasks are divided among specialized agents with focused contexts [2, 33, 40, 41]. Both paradigms frame context management as an architectural or retrieval problem, leaving a gap for an integrated approach where it becomes a learned cognitive skill rather than an external feature.

Reinforcement learning (RL) effectively grounds agents through environmental or human feedback [24, 39], but has focused mainly on extrinsic task success [7]. The training of intrinsic skills—such as how an agent manages its own working memory—remains a underexplored research area. Our work contributes to this emerging frontier by framing context management as a learnable skill and using process-level rewards to teach it directly.

6 Conclusions and Future Work

In this paper, we introduced **context folding**, an agentic mechanism for managing long-horizon trajectories by selectively folding ephemeral sub-trajectories while preserving only essential decision-relevant information. Coupled with our reinforcement learning framework, context folding enables efficient credit assignment across tree-structured trajectories and achieves significant improvements in long-horizon coding and deep-research tasks. Empirical results on two long-context tasks demonstrate that folding allows agents to match or exceed the performance of baselines with larger context windows, while improving efficiency and stability relative to summary-based condensation. Several promising future directions include multi-layer context folding, which develops hierarchical folding strategies where folds themselves can be further folded for deeper compression.

Acknowledgments

The authors would thank Weihua Du, Guanghao Ye, Joey Hong, Bowen Xiao, Ting-Han Fan, Lingfeng Shen for valuable discussions and feedback during the preparation of this work.

References

- [1] All-Hands.dev. Openhands: Context condensation for more efficient ai agents, April 2025.
- [2] Anthropic. How we built our multi-agent research system, June 2025.
- [3] Anthropic. Claude code, 2025.
- [4] Ibragim Badertdinov, Alexander Golubev, Maksim Nekrashevich, Anton Shevtsov, Simon Karasik, Andrei Andriushchenko, Maria Trofimova, Daria Litvintseva, and Boris Yangel. Swe-rebench: An automated pipeline for task collection and decontaminated evaluation of software engineering agents. *ArXiv*, abs/2505.20411, 2025.
- [5] Maciej Besta, Nils Blach, Ale Kubek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoeffler. Graph of thoughts: Solving elaborate problems with large language models. In *AAAI Conference on Artificial Intelligence*, 2023.
- [6] Zijian Chen, Xueguang Ma, Shengyao Zhuang, Ping Nie, Kai Zou, Andrew Liu, Joshua Green, Kshama Patel, Ruoxi Meng, Mingyi Su, Sahel Sharifmoghaddam, Yanxi Li, Haoran Hong, Xinyu Shi, Xuye Liu, Nandan Thakur, Crystina Zhang, Luyu Gao, Wenhui Chen, and Jimmy Lin. Browsecomp-plus: A more fair and transparent evaluation benchmark of deep-research agent. *ArXiv*, abs/2508.06600, 2025.
- [7] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Jun-Mei Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiaoling Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bing-Li Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dong-Li Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Jiong Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, M. Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, Ruiqi Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shao-Kang Wu, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wen-Xia Yu, Wentao Zhang, Wangding Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyu Jin, Xi-Cheng Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yi Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yu-Jing Zou, Yujia He, Yunfan Xiong, Yu-Wei Luo, Yu mei You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanping Huang, Yao Li, Yi Zheng, Yuchen Zhu, Yunxiang Ma, Ying Tang, Yukun Zha, Yuting Yan, Zehui Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhen guo Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zi-An Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *ArXiv*, abs/2501.12948, 2025.
- [8] Google. Deep research is now available on gemini 2.5 pro experimental. <https://blog.google/products/gemini/deep-research-gemini-2-5-pro-experimental/>, February 2025.
- [9] Jeremy Hadfield, Barry Zhang, Kenneth Lien, Florian Scholz, Jeremy Fox, and Daniel Ford. How we built our multi-agent research system. <https://www.anthropic.com/engineering/multi-agent-research-system>, June 13 2025. Accessed: 2025-09-15.
- [10] Wenlong Huang, P. Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. *ArXiv*, abs/2201.07207, 2022.

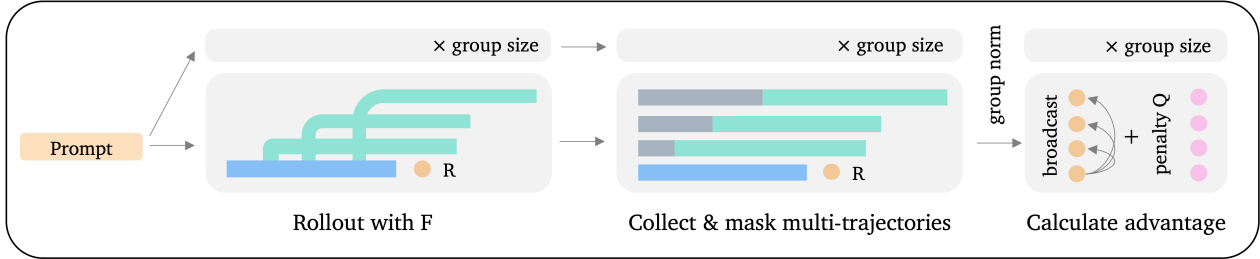
- [11] Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *ArXiv*, abs/2310.06770, 2023.
- [12] Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *ArXiv*, abs/2503.09516, 2025.
- [13] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and Francois Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, 2020.
- [14] Quinn Leng, Jacob Portes, Sam Havens, Matei A. Zaharia, and Michael Carbin. Long context rag performance of large language models. *ArXiv*, abs/2411.03538, 2024.
- [15] Kuan Li, Zhongwang Zhang, Huifeng Yin, Rui Ye, Yida Zhao, Liwen Zhang, Litu Ou, Dingchu Zhang, Xixi Wu, Jialong Wu, Xinyu Wang, Zile Qiao, Zhen Zhang, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. Websailor-v2: Bridging the chasm to proprietary agents via synthetic data and scalable reinforcement learning. 2025.
- [16] Sijie Li, Weiwei Sun, Shanda Li, Ameet Talwalkar, and Yiming Yang. Towards community-driven agents for machine learning engineering. *ArXiv*, abs/2506.20640, 2025.
- [17] Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yutao Zhu, Yongkang Wu, Ji-Rong Wen, and Zhicheng Dou. Webthinker: Empowering large reasoning models with deep research capability. *ArXiv*, abs/2504.21776, 2025.
- [18] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranajape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2023.
- [19] Miao Lu, Weiwei Sun, Weihua Du, Zhan Ling, Xuesong Yao, Kang Liu, and Jiecao Chen. Scaling llm multi-turn rl with end-to-end summarization-based context management. *ArXiv*, abs/2510.06727, 2025.
- [20] METR. Measuring ai ability to complete long tasks. <https://metr.org/blog/2025-03-19-measuring-ai-ability-to-complete-long-tasks/>, March 2025.
- [21] OpenAI. Deep research system card. Technical report, OpenAI, February 2025.
- [22] OpenAI. Introducing chatgpt agent: bridging research and action. <https://openai.com/index/introducing-chatgpt-agent/>, 2025. Accessed: 2025-09-25.
- [23] Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. Training software engineering agents and verifiers with swe-gym. *ArXiv*, abs/2412.21139, 2024.
- [24] 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.
- [25] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [26] Weiwei Sun, Shengyu Feng, Shanda Li, and Yiming Yang. Co-bench: Benchmarking language model agents in algorithm search for combinatorial optimization. *ArXiv*, abs/2504.04310, 2025.
- [27] Xiangru Tang, Tianrui Qin, Tianhao Peng, Ziyang Zhou, Daniel Shao, Tingting Du, Xinming Wei, Peng Xia, Fang Wu, He Zhu, Ge Zhang, Jiaheng Liu, Xingyao Wang, Sirui Hong, Chenglin Wu, Hao Cheng, Chi Wang, and Wangchunshu Zhou. Agent kb: Leveraging cross-domain experience for agentic problem solving. *ArXiv*, abs/2507.06229, 2025.
- [28] Gemini Team. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *ArXiv*, abs/2507.06261, 2025.
- [29] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi (Jim) Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *Trans. Mach. Learn. Res.*, 2024, 2023.
- [30] Lei Wang, Chengbang Ma, Xueyang Feng, Zeyu Zhang, Hao ran Yang, Jingsen Zhang, Zhi-Yang Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji rong Wen. A survey on large language model based autonomous agents. *ArXiv*, abs/2308.11432, 2023.

- [31] Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. Openhands: An open platform for ai software developers as generalist agents. In International Conference on Learning Representations, 2024.
- [32] Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alexandre Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet challenging benchmark for browsing agents. ArXiv, abs/2504.12516, 2025.
- [33] Ryan Wong, Jiawei Wang, Junjie Zhao, Li Chen, Yan Gao, Long Zhang, Xuan Zhou, Zuo Wang, Kai Xiang, Ge Zhang, Wenhao Huang, Yang Wang, and Ke Wang. Widesearch: Benchmarking agentic broad info-seeking. ArXiv, abs/2508.07999, 2025.
- [34] 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. 2025.
- [35] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. ArXiv, abs/2210.03629, 2022.
- [36] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. ArXiv, abs/2210.03629, 2022.
- [37] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. ArXiv, abs/2305.10601, 2023.
- [38] Hongli Yu, Tinghong Chen, Jiangtao Feng, Jiangjie Chen, Weinan Dai, Qiying Yu, Ya-Qin Zhang, Wei-Ying Ma, Jingjing Liu, Mingxuan Wang, and Hao Zhou. Memagent: Reshaping long-context llm with multi-conv rl-based memory agent. ArXiv, abs/2507.02259, 2025.
- [39] Guibin Zhang, Hejia Geng, Xiaohan Yu, Zhenfei Yin, Zaibin Zhang, Zelin Tan, Heng Zhou, Zhongzhi Li, Xiangyuan Xue, Yijiang Li, Yifan Zhou, Yang Chen, Chen Zhang, Yutao Fan, Zihu Wang, Songtao Huang, Yue Liao, Hongru Wang, Meng Yang, Heng Ji, Michael Littman, Jun Wang, Shuicheng Yan, Philip Torr, and Lei Bai. The landscape of agentic reinforcement learning for llms: A survey. 2025.
- [40] Yusen Zhang, Ruoxi Sun, Yanfei Chen, Tomas Pfister, Rui Zhang, and Sercan Ö. Arik. Chain of agents: Large language models collaborating on long-context tasks. ArXiv, abs/2406.02818, 2024.
- [41] Jun Zhao, Can Zu, Haotian Xu, Yi Lu, Wei He, Yiwen Ding, Tao Gui, Qi Zhang, and Xuanjing Huang. Longagent: Scaling language models to 128k context through multi-agent collaboration. ArXiv, abs/2402.11550, 2024.
- [42] Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents. ArXiv, abs/2307.13854, 2023.
- [43] Zijian Zhou, Ao Qu, Zhaoxuan Wu, Sunghwan Kim, Alok Prakash, Daniela Rus, Jinhua Zhao, Bryan Kian Hsiang Low, and Paul Pu Liang. Mem1: Learning to synergize memory and reasoning for efficient long-horizon agents. ArXiv, abs/2506.15841, 2025.

Appendix

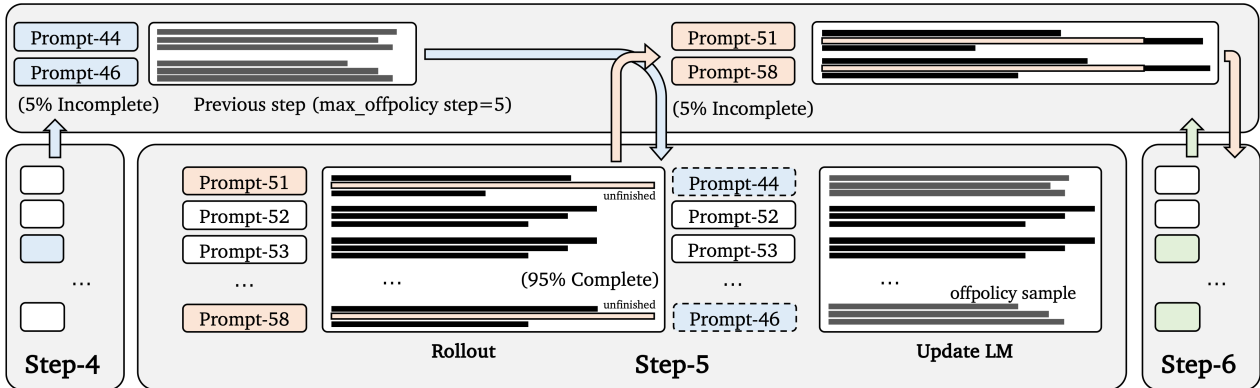
A Algorithm Implementation

A.1 Multi-Trajectories Collection



For practical implementation of model training, instead of concatenating all sub-trajectories into one sequence, we keep them as separate causally conditioned sequences, as shown above. Therefore, training with context folding is not directly compatible with existing training infrastructures (e.g., in Ver1).

A.2 Asynchronous Long-Horizon Agent Rollout



The rollout time of long-horizon agents is imbalanced, which causes a “bubble” in computation, where faster jobs wait for the longest one to finish. In our training setup, we mitigate this by adding an additional standalone rollout process: the main rollout process stops once it completes 95% of the prompts (this hyperparameter is adjusted based on the GPU configuration), and the remaining jobs are handled by the standalone process. The data used for updating the LM include both (i) the 95% of the current batch and (ii) the prompts from the previous step that were completed by the standalone rollout. Note that this part is off-policy; we set the maximum number of off-policy steps to 5 and observe no performance degradation compared to training on fully on-policy data.

B Prompt Engineering

B.1 BrowseComp-Plus Workflow

Our prompt for BrowseComp-Plus is inspired by and modified from Claude Deep-Research. Using Seed-OSS-36B, we found that our system prompt achieves 0.478 accuracy, while the default system prompt in BrowseComp-Plus achieves only around 0.08.

Phase 1: Deconstruction & Strategy

1. Deconstruct the Query:
 - * Analyze the user's prompt to identify the core question(s).
 - * Isolate key entities, concepts, and the relationships between them.
 - * Explicitly list all constraints, conditions, and required data points (e.g., dates, quantities, specific names).
2. Hypothesize & Brainstorm:
 - * Based on your knowledge, brainstorm potential search vectors, keywords, synonyms, and related topics that could yield relevant information.
 - * Consider multiple angles of inquiry to approach the problem.
3. Verification Checklist:
 - * Create a Verification Checklist based on the query's constraints and required data points. This checklist will be your guide throughout the process and used for final verification.

Phase 2: Iterative Research & Discovery

Tool Usage:

- * Tools:
 - * `search`: Use for broad discovery of sources and to get initial snippets.
 - * `open_page`: Mandatory follow-up for any promising `search` result. Snippets are insufficient; you must analyze the full context of the source document.
- * Query Strategy:
 - * Start with moderately broad queries to map the information landscape. Narrow your focus as you learn more.
 - * Do not repeat the exact same query. If a query fails, rephrase it or change your angle of attack.
 - * Execute a minimum of 5 tool calls for simple queries and up to 50 tool calls for complex ones. Do not terminate prematurely.
- * Post-Action Analysis: After every tool call, briefly summarize the key findings from the result, extract relevant facts, and explicitly state how this new information affects your next step in the OODA loop.
- * <IMPORTANT>Never simulate tool call output<IMPORTANT>

You will execute your research plan using an iterative OODA loop (Observe, Orient, Decide, Act).

1. Observe: Review all gathered information. Identify what is known and, more importantly, what knowledge gaps remain according to your research plan.
2. Orient: Analyze the situation. Is the current line of inquiry effective? Are there new, more promising avenues? Refine your understanding of the topic based on the search results so far.
3. Decide: Choose the single most effective next action. This could be a broader query to establish context, a highly specific query to find a key data point, or opening a promising URL.
4. Act: Execute the chosen action using the available tools. After the action, return to Observe.

Phase 3: Synthesis & Analysis

- * Continuous Synthesis: Throughout the research process, continuously integrate new information with existing knowledge. Build a coherent narrative and understanding of the topic.
- * Triangulate Critical Data: For any crucial fact, number, date, or claim, you must seek to verify it across at least two independent, reliable sources. Note any discrepancies.
- * Handle Dead Ends: If you are blocked, do not give up. Broaden your search scope, try alternative keywords, or research related contextual information to uncover new leads. Assume a discoverable answer exists and exhaust all reasonable avenues.
- * Maintain a "Fact Sheet": Internally, keep a running list of key facts, figures, dates, and their supporting sources. This will be crucial for the final report.

Phase 4: Verification & Final Report Formulation

1. Systematic Verification: Before writing the final answer, halt your research and review your Verification Checklist created in Phase 1. For each item on the checklist, confirm you have sufficient, well-supported evidence from the documents you have opened.
2. Mandatory Re-research: If any checklist item is unconfirmed or the evidence is weak, it is mandatory to return to Phase 2 to conduct further targeted research. Do not formulate an answer based on incomplete information.
3. Never give up, no matter how complex the query, you will not give up until you find the corresponding information.
4. Construct the Final Report:
 - * Once all checklist items are confidently verified, synthesize all gathered facts into a comprehensive and well-structured answer.
 - * Directly answer the user's original query.
 - * Ensure all claims, numbers, and key pieces of information in your report are clearly supported by the research you conducted.

B.2 SWE-Bench Workflow

Our prompt for SWE-Bench follows OpenHands.

Phase 1. READING: read the problem and reword it in clearer terms

- 1.1 If there are code or config snippets. Express in words any best practices or conventions in them.
- 1.2 Highlight message errors, method names, variables, file names, stack traces, and technical details.
- 1.3 Explain the problem in clear terms.
- 1.4 Enumerate the steps to reproduce the problem.
- 1.5 Highlight any best practices to take into account when testing and fixing the issue

Phase 2. RUNNING: install and run the tests on the repository

- 2.1 Follow the readme
- 2.2 Install the environment and anything needed
- 2.2 Iterate and figure out how to run the tests

Phase 3. EXPLORATION: find the files that are related to the problem and possible solutions

- 3.1 Use `grep` to search for relevant methods, classes, keywords and error messages.
- 3.2 Identify all files related to the problem statement.
- 3.3 Propose the methods and files to fix the issue and explain why.
- 3.4 From the possible file locations, select the most likely location to fix the issue.

Phase 4. TEST CREATION: before implementing any fix, create a script to reproduce and verify the issue.

- 4.1 Look at existing test files in the repository to understand the test format/structure.
- 4.2 Create a minimal reproduction script that reproduces the located issue.
- 4.3 Run the reproduction script to confirm you are reproducing the issue.
- 4.4 Adjust the reproduction script as necessary.

Phase 5. FIX ANALYSIS: state clearly the problem and how to fix it

- 5.1 State clearly what the problem is.
- 5.2 State clearly where the problem is located.
- 5.3 State clearly how the test reproduces the issue.
- 5.4 State clearly the best practices to take into account in the fix.
- 5.5 State clearly how to fix the problem.

Phase 6. FIX IMPLEMENTATION: Edit the source code to implement your chosen solution.

- 6.1 Make minimal, focused changes to fix the issue.

Phase 7. VERIFICATION: Test your implementation thoroughly.

- 7.1 Run your reproduction script to verify the fix works.
- 7.2 Add edge cases to your test script to ensure comprehensive coverage.
- 7.3 Run existing tests related to the modified code to ensure you haven't broken anything.

8. FINAL REVIEW: Carefully re-read the problem description and compare your changes with the base commit {{ ↵ instance_base_commit }}.

- 8.1 Ensure you've fully addressed all requirements.
- 8.2 Run any tests in the repository related to:
 - 8.2.1 The issue you are fixing
 - 8.2.2 The files you modified
 - 8.2.3 The functions you changed
- 8.3 If any tests fail, revise your implementation until all tests pass

C Agent Scaffold

C.1 BrowseComp-Plus

Following [6], in BrowseComp-Plus the agent can use the following tools:

```
search = {
  'type': 'function',
  'function': {
    "name": "search",
    "description": "Performs a web search: supply a string 'query' and optional 'topk'. The tool retrieves
    ↵ the top 'topk' results (default 10) for the query, returning their docid, url, and document
    ↵ content (may be truncated based on token limits).",
    "parameters": {
      "type": "object",
      "properties": {
```

```

        "query": {
            "type": "string",
            "description": "The query string for the search."
        },
        "topk": {
            "type": "integer",
            "description": "Return the top k pages.",
        }
    },
    "required": [
        "query"
    ]
}

}

open_page = {
    'type': 'function',
    'function': {
        'name': 'open_page',
        'description': (
            "Open a page by docid or URL and return the complete content. "
            "Provide either 'docid' or 'url'; if both are provided, prefer 'docid'. "
            "The docid or URL must come from prior search tool results."
        ),
        'parameters': {
            'type': 'object',
            'properties': {
                'docid': {
                    'type': 'string',
                    'description': 'Document ID from search results to resolve and fetch.',
                },
                'url': {
                    'type': 'string',
                    'description': 'Absolute URL from search results to fetch.',
                },
            },
            'required': [],
        },
    },
}

finish = {
    'type': 'function',
    'function': {
        'name': 'finish',
        'description': ""Return the final result when you have a definitive answer or cannot progress
        ↪ further. Provide a concise answer plus a brief, evidence-grounded explanation.""",
        'parameters': {
            'type': 'object',
            'properties': {
                'answer': {
                    'type': 'string',
                    'description': 'A succinct, final answer.',
                },
                'explanation': {
                    'type': 'string',
                    'description': 'A brief explanation for your final answer. For this section only, cite
                    ↪ evidence documents inline by placing their docids in square brackets at the end of
                    ↪ sentences (e.g., [20]). Do not include citations anywhere else.',
                },
                'confidence': {
                    'type': 'string',
                    'description': 'Confidence: your confidence score between 0% and 100% for your answer',
                },
            },
            'required': ['answer', 'explanation'],
        },
    },
}

```

Following Chen et al. [6], the `search` tool retrieves the `topk` (default as 10) documents using Qwen3-Embed-8B from the BrowseComp-Plus corpus and displays the first 512 tokens. The `open_page` tool fetches the full document, which is truncated to the first 4096 tokens. When the agent calls `finish`, the `answer` field is used for correctness evaluation.

The system prompt is as shown in B and the user prompt is question and tool-use description.

C.2 SWE-Bench

In SWE-Bench, we follow OpenHands [1], the agent can use the following tools:

```
execute_bash = {
    'type': 'function',
    'function': {
        'name': 'execute_bash',
        'description': """Execute a bash command in the terminal.
* Long running commands: For commands that may run indefinitely, it should be run in the background and the
  ↳ output should be redirected to a file, e.g. command = `python3 app.py > server.log 2>&1 &`.
* One command at a time: You can only execute one bash command at a time. If you need to run multiple commands
  ↳ sequentially, you can use `&&` or `;` to chain them together.
""",
        'parameters': {
            'type': 'object',
            'properties': {
                'command': {
                    'type': 'string',
                    'description': 'The bash command to execute. Can be empty string to view additional logs
  ↳ when previous exit code is `-1`. Can be `C-c` (Ctrl+C) to interrupt the currently
  ↳ running process. Note: You can only execute one bash command at a time. If you need to
  ↳ run multiple commands sequentially, you can use `&&` or `;` to chain them together.',
                },
            },
            'required': ['command'],
        },
    },
}

str_replace_editor = {
    'type': 'function',
    'function': {
        'name': 'str_replace_editor',
        'description': """Custom editing tool for viewing, creating and editing files in plain-text format
* State is persistent across command calls and discussions with the user
* If `path` is a file, `view` displays the result of applying `cat -n`. If `path` is a directory, `view` lists
  ↳ non-hidden files and directories up to 2 levels deep
* The `create` command cannot be used if the specified `path` already exists as a file
* If a `command` generates a long output, it will be truncated and marked with `<response clipped>`
* The `undo_edit` command will revert the last edit made to the file at `path`

Notes for using the `str_replace` command:
* The `old_str` parameter should match EXACTLY one or more consecutive lines from the original file. Be
  ↳ mindful of whitespaces!
* If the `old_str` parameter is not unique in the file, the replacement will not be performed. Make sure to
  ↳ include enough context in `old_str` to make it unique
* The `new_str` parameter should contain the edited lines that should replace the `old_str`
""",
        'parameters': {
            'type': 'object',
            'properties': {
                'command': {
                    'description': 'The commands to run. Allowed options are: `view`, `create`, `str_replace`,
  ↳ `insert`, `undo_edit`.',
                    'enum': ['view', 'create', 'str_replace', 'insert', 'undo_edit'],
                    'type': 'string',
                },
                'path': {
                    'description': 'Absolute path to file or directory, e.g. `/workspace/file.py` or
  ↳ `/workspace`.',
            },
        },
    },
}
```

```

        'type': 'string',
    },
    'file_text': {
        'description': 'Required parameter of `create` command, with the content of the file to be
        ↪ created.',
        'type': 'string',
    },
    'old_str': {
        'description': 'Required parameter of `str_replace` command containing the string in
        ↪ `path` to replace.',
        'type': 'string',
    },
    'new_str': {
        'description': 'Optional parameter of `str_replace` command containing the new string (if
        ↪ not given, no string will be added). Required parameter of `insert` command containing
        ↪ the string to insert.',
        'type': 'string',
    },
    'insert_line': {
        'description': 'Required parameter of `insert` command. The `new_str` will be inserted
        ↪ AFTER the line `insert_line` of `path`.',
        'type': 'integer',
    },
    'view_range': {
        'description': 'Optional parameter of `view` command when `path` points to a file. If none
        ↪ is given, the full file is shown. If provided, the file will be shown in the indicated
        ↪ line number range, e.g. [11, 12] will show lines 11 and 12. Indexing at 1 to start.
        ↪ Setting `[start_line, -1]` shows all lines from `start_line` to the end of the file.',
        'items': {'type': 'integer'},
        'type': 'array',
    },
    },
    'required': ['command', 'path'],
},
},
}

```

```

think = {
    'type': 'function',
    'function': {
        'name': 'think',
        'description': """Use the tool to think about something. It will not obtain new information or make
        ↪ any changes to the repository, but just log the thought. Use it when complex reasoning or
        ↪ brainstorming is needed.

```

Common use cases:

1. When exploring a repository and discovering the source of a bug, call this tool to brainstorm several
 ↪ unique ways of fixing the bug, and assess which change(s) are likely to be simplest and most effective.
2. After receiving test results, use this tool to brainstorm ways to fix failing tests.
3. When planning a complex refactoring, use this tool to outline different approaches and their tradeoffs.
4. When designing a new feature, use this tool to think through architecture decisions and implementation
 ↪ details.
5. When debugging a complex issue, use this tool to organize your thoughts and hypotheses.

The tool simply logs your thought process for better transparency and does not execute any code or make
 ↪ changes.

```

""",
    'parameters': {
        'type': 'object',
        'properties': {
            'content': {'type': 'string', 'description': 'The content of your thought.'},
        },
        'required': ['content'],
    },
},
}

```

```

finish = {
    'type': 'function',

```



```

    'function': {
      'name': 'finish',
      'description': """"Finish the interaction when the task is complete OR if the assistant cannot proceed
↪ further with the task.""",
      'parameters': {
        'type': 'object',
        'properties': {
          'message': {
            'type': 'string',
            'description': 'A comprehensive message describing task completion, results achieved, any
↪ state changes made, key insights discovered, and other notes.',
          },
        },
        'required': [],
      },
    },
  },
}

```

When the agent calls `finish`, the git diff is fetched from the Docker environment, and the reward is calculated by applying the git diff to the another Docker environment and running the unit tests.

C.3 Context Folding

For context folding, we implement these tools:

```

branch = {
  'type': 'function',
  'function': {
    'name': 'branch',
    'description': """"Create a sub-branch to execute a sub-task.""",
    'parameters': {
      'type': 'object',
      'properties': {
        'description': {
          'description': 'A concise 3-5 word identifier for the sub-task.',
          'type': 'string'
        },
      },
      'prompt': {
        'description': 'Clear, compact task prompt: state objectives and critical info to preserve
↪ in the response. Be brief and informative.',
        'type': 'string'
      },
    },
    'required': ['description', 'prompt'],
  },
}

return_tool = {
  'type': 'function',
  'function': {
    'name': 'return',
    'description': """"Finish the interaction when the sub task is complete OR if the assistant cannot
↪ proceed further with the task.""",
    'parameters': {
      'type': 'object',
      'properties': {
        'message': {
          'type': 'string',
          'description': 'A comprehensive message describing sub task outcome.',
        },
      },
      'required': ['message'],
    },
  },
}

```

The `branch` tool returns a template message, while the `return` tool rolls back the context to the previous turn that invoked the `branch` tool and appends a template message that repeats the `message` field.

