

RESUM: UNLOCKING LONG-HORIZON SEARCH INTELLIGENCE VIA CONTEXT SUMMARIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Model (LLM)-based web agents demonstrate strong performance on knowledge-intensive tasks but are hindered by context window limitations in paradigms like ReAct. Complex queries involving multiple entities, intertwined relationships, and high uncertainty demand extensive search cycles that rapidly exhaust context budgets before reaching solutions. To overcome this challenge, we introduce ReSum, a novel paradigm that enables indefinite exploration through periodic context summarization. ReSum converts growing interaction histories into compact reasoning states, maintaining awareness of prior discoveries while bypassing context constraints. For paradigm adaptation, we propose ReSum-GRPO, integrating GRPO with segmented trajectory training and advantage broadcasting to familiarize agents with summary-conditioned reasoning. Extensive experiments on web agents across three benchmarks demonstrate that ReSum delivers an average absolute improvement of 4.5% over ReAct, with further gains of 8.2% following ReSum-GRPO training. Notably, with only 1K training samples, the ReSum-GRPO-trained 30B model achieves 33.3% Pass@1 on BrowseComp-zh and 18.3% on BrowseComp-en, [showing competitive performance with leading open-source web agents](#).

1 INTRODUCTION

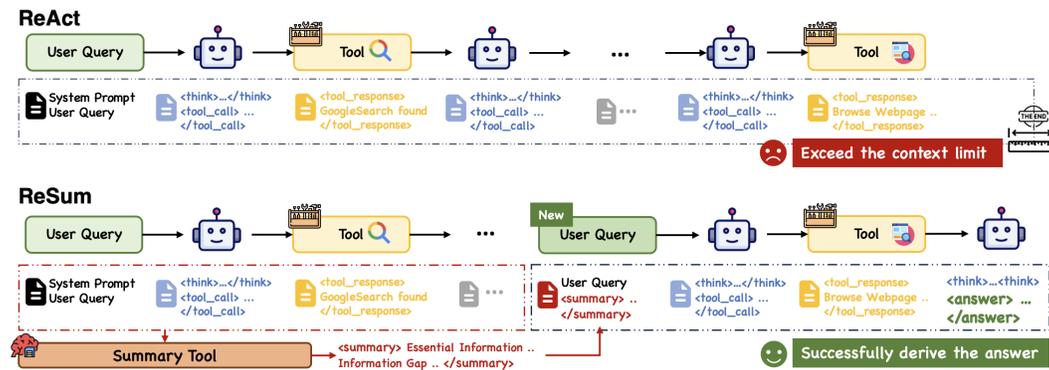


Figure 1: **Comparison between ReAct and ReSum paradigms.** Appending every observation, thought, and action in ReAct exhausts the context budget before multi-turn exploration completes. In contrast, ReSum periodically invokes a summary tool to condense history and resumes reasoning from the compressed summary, enabling indefinite exploration.

Recent advances in Large Language Model (LLM)-based agents have demonstrated strong performance on complex, knowledge-intensive tasks (Yao et al., 2023; Wang et al., 2024; Jin et al., 2025; Xi et al., 2025; Gao et al., 2025). Among these, **web agents** are particularly critical: they actively search and browse the open web, extract and ground facts from diverse sources, and synthesize answers that are both user-specific and up-to-date (Wu et al., 2025a; Li et al., 2025a; Tao et al., 2025).

Achieving reliable and comprehensive answers for complex queries is nontrivial. Consider this example: “A painter, whose father died of heart disease, had an elder sister and five children with

054 *his wife. Later, his marriage broke down and he had three more relationships. What is the name*
 055 *of the literary work based on this person?"* This query exemplifies the challenges of complex web
 056 search: it involves **multiple entities**, **intertwined relationships**, and fragmentary information with
 057 **high uncertainty**. Such questions cannot be resolved with a few search calls. Instead, they require
 058 extended cycles of targeted querying, browsing, extraction, and cross-verification to progressively
 059 reduce uncertainty and assemble a complete and grounded evidence chain (Gao et al., 2025).

060 However, such long-horizon exploration faces a fundamental barrier: context constraints. Most LLMs
 061 have limited context windows, e.g., 32k tokens (Yang et al., 2024; Team, 2025b; Jiang et al., 2023),
 062 and the popular ReAct paradigm (Yao et al., 2023), which appends every observation, thought, and
 063 action to the dialogue history, quickly exhausts this budget before answers are found (illustrated in
 064 **Figure 1**). As trajectories lengthen, accumulated context forces premature termination, blocking
 065 solutions that require extensive exploration.

066 To address this challenge, we introduce **ReSum**, a novel paradigm enabling indefinite exploration
 067 through context summarization. The core insight is to convert growing interaction histories into
 068 compact reasoning states before context limits are exceeded. Rather than appending every interaction,
 069 ReSum *periodically* compresses conversations into structured summaries and resumes exploration
 070 from these states, allowing agents to maintain awareness of prior discoveries without context con-
 071 straints. **Unlike prior context management approaches that require architectural modifications and**
 072 **end-to-end training of new agents** (Zhou et al., 2025b; Yu et al., 2025a), ReSum introduces a **paradigm-**
 073 **matic enhancement through minimal modifications to ReAct**, maintaining simplicity and ensuring
 074 compatibility with pre-existing agents.

075 ReSum utilizes an off-the-shelf LLM as the summary tool, but generic LLMs, particularly smaller
 076 models, often struggle with effectively summarizing conversations in web search contexts. Therefore,
 077 we specialize the summarization capability by fine-tuning Qwen3-30B-A3B-Thinking (Team, 2025b)
 078 using `<Conversation, Summary>` pairs collected from powerful open-source models (Guo et al.,
 079 2025; OpenAI, 2025a), resulting in ReSumTool-30B. Unlike conventional summarization tools,
 080 ReSumTool-30B is specifically trained to **extract key clues and evidence** from lengthy interactions,
 081 **identify information gaps**, and highlight next-step directions. This specialization makes it uniquely
 082 suited for web search tasks, combining lightweight deployment with task-specific enhancements.
 083 Extensive evaluations demonstrate that ReSumTool-30B outperforms larger models like as Qwen3-
 084 235B (Team, 2025b) and DeepSeek-R1-671B (Guo et al., 2025) in summarization quality.

085 Finally, to enable agents to master the ReSum paradigm, we employ reinforcement learning (RL)
 086 through our tailored ReSum-GRPO algorithm. Unlike supervised fine-tuning, which demands costly
 087 expert-level ReSum trajectory data and risks overwriting an agent’s existing skills, RL allows agents
 088 to adapt to the paradigm through self-evolution without compromising their inherent reasoning
 089 capabilities. Specifically, ReSum-GRPO follows standard GRPO procedures (Shao et al., 2024) with
 090 modifications for long trajectories: when approaching context limits, agents invoke ReSumTool-
 091 30B to compress the conversation and continue from the summary state, naturally segmenting the
 092 complete trajectory into multiple parts. Each segment becomes an individual training episode, and we
 093 broadcast the trajectory-level advantage to *all* segments within the same trajectory. This mechanism
 094 encourages agents to both reason effectively from compressed states and collect information that
 095 produces high-quality summaries. Our main contributions are summarized as follows:

- 096 • **ReSum: A paradigm for long-horizon web search** We identify the fundamental context
 097 limitation of the ReAct paradigm and propose ReSum, which periodically compresses
 098 conversation history into compact summary, enabling agents to indefinitely explore. **ReSum**
 099 **provides a plug-and-play solution to existing ReAct agents, contrasting with previous**
 100 **architecture-modifying approaches that necessitate end-to-end training of new agents.**
- 101 • **ReSumTool-30B: A specialized summary model** To enable goal-oriented conversation
 102 summarization, we develop ReSumTool-30B through targeted training **with rigorous quality**
 103 **control**, designed specifically to extract key evidence from lengthy interactions.
- 104 • **ReSum-GRPO: Reinforcement learning for paradigm adaptation** We design ReSum-
 105 GRPO to familiarize agents with summary-based reasoning by segmenting long trajectories
 106 and broadcasting trajectory-level advantages across all segments. Experiments across web
 107 agents on three challenging benchmarks show average improvements of 4.5% for ReSum
 compared to ReAct, with further improvements of 8.2% after ReSum-GRPO training.

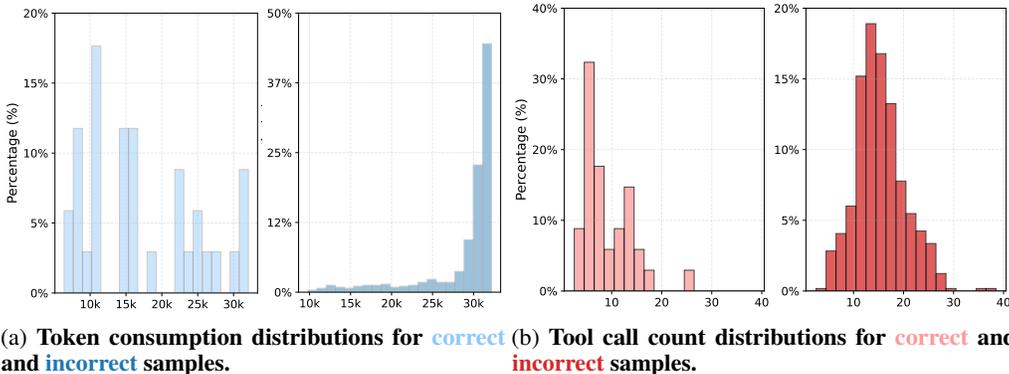


Figure 2: **Context limits in ReAct constrain exploration.** Using open-sourced WebSailor-7B (Li et al., 2025a) on the BrowseComp-en (Wei et al., 2025), we compare the distributions of token consumption and tool call counts between correctly solved and failed trajectories. Failed cases use far more tool calls and tokens, suggesting trajectories are frequently truncated before resolution.

2 PRELIMINARY

Before introducing ReSum, we first review the popular ReAct paradigm to facilitate comprehension and highlight the fundamental challenges that motivate our research.

ReAct (Yao et al., 2023) is a widely adopted agentic workflow (Li et al., 2025a; Wu et al., 2025a; Li et al., 2025c) where agents perform iterative cycles of Thought, Action, and Observation. Specifically, in each iteration, the LLM generates a reasoning step (Thought) based on existing context, executes a parsable tool call (Action), and receives feedback from the environment (Observation). In web search contexts, the action space typically consists of search queries, webpage browsing, or generating the final answer. The iteration terminates when the agent produces a final answer. A complete trajectory with T iterations can be formally defined as:

$$\mathcal{H}_T = (q, \tau_1, a_1, o_1, \dots, \tau_{T-1}, a_{T-1}, o_{T-1}, \tau_T, a_T),$$

where q is the question, and τ_i , a_i , and o_i represent the thought, action, and observation at the i -th round, respectively. At step t , both the thought τ_t and action a_t are sampled from a policy model π_θ conditioned on all previous context as $(\tau_t, a_t) \sim \pi_\theta(\cdot | \mathcal{H}_{t-1})$.

For complex web search tasks with highly ambiguous entities and relationships, agents must perform extensive tool interactions to gather disparate evidence and converge on a solution. However, continuously appending the full interaction history quickly exhausts modern LLMs’ context windows before difficult cases can be resolved. To illustrate this limitation, we analyze the WebSailor-7B agent’s behavior (Li et al., 2025a) on the challenging BrowseComp-en benchmark (Wei et al., 2025). As shown in Figures 2a and 2b, most solved cases complete within 10 tool calls, whereas failed cases typically exceed 10, and often 20, resulting in sharply increased token usage that surpasses the 32k limit.

Following existing web agent designs (Li et al., 2025a; Gao et al., 2025; Li et al., 2025b), we implement two essential tools for web exploration: **Search** queries the Google Search engine, accepting multiple queries simultaneously and returning the top-10 results per query, and **Visit** browses specific web pages by URL using Jina (Jina.ai, 2025) and extracts goal-specific evidence using Qwen2.5-72B-Instruct (Yang et al., 2024). Due to space constraints, the discussion of existing research on web agents and context management for agents has been moved to **Appendix A**.

3 METHODOLOGY

In this section, we introduce the ReSum paradigm, the development of ReSumTool-30B, and the ReSum-GRPO algorithm designed to facilitate paradigm adaptation.

3.1 RESUM PARADIGM

Trajectory Initialization: The trajectory begins with a user query q , initializing $\mathcal{H}_0 = (q)$. Following ReAct, the agent alternates between internal reasoning and tool use: at the t -th round, it generates a reasoning step within `<think>` `</think>` tokens and issues a tool call within `<tool_call>` `</tool_call>` tokens, expressed as $(\tau_t, a_t) \sim \pi_\theta(\cdot \mid \mathcal{H}_{t-1})$. The system parses the tool call arguments and executes the corresponding tool, returning results within `<tool_response>` `</tool_response>` tokens as $o_t = \mathcal{R}(a_t)$, where \mathcal{R} represents the tool environment. The history is then updated by concatenation as follows:

$$\mathcal{H}_t = \mathcal{H}_{t-1} \circ (\tau_t, a_t, o_t).$$

In the initial rounds, ReSum mirrors ReAct by iteratively building $\mathcal{H}_t = (q, \tau_1, a_1, o_1, \dots, \tau_t, a_t, o_t)$.

Context Summarization: When a compression trigger is activated, a summary tool π_{sum} is invoked to summarize the accumulated history as:

$$s \sim \pi_{\text{sum}}(\cdot \mid \mathcal{H}_t),$$

where s is a goal-oriented `<summary>` `</summary>` that consolidates verified evidence and explicitly lists information gaps (prompt provided in Appendix C). Then, we form a compressed state $q' = (q, s)$ and reset the working history to:

$$\mathcal{H}_t \leftarrow (q').$$

Triggers for summarization can be *systematic*, e.g., exceeding a token budget or reaching a round limit, or *agent-initiated*, where the policy model decides to summarize for effective context management.

Trajectory Termination: Through periodic summarization, ReSum dynamically maintains the context within the model’s window while preserving essential evidence. The agent continues gathering evidence and, once sufficient information is accumulated, produces a synthesized answer within `<answer>` `</answer>` tokens. Although ReSum theoretically allows unbounded exploration, practical deployments impose resource budgets, e.g., limiting the number of tool calls. Trajectories that exceed these limits are terminated and marked as failures.

The complete ReSum workflow is detailed in **Algorithm 1** in Appendix. Unlike ReAct, which accumulates all interactions, ReSum transforms lengthy interaction histories into compact, restartable reasoning states. This approach distills key evidence and highlights actionable next steps, enabling multi-turn exploration while adhering to token budget constraints. Furthermore, by minimizing modifications to ReAct, ReSum preserves its simplicity, efficiency, and functional compatibility with existing ReAct-trained agents.

3.2 SUMMARY TOOL SPECIFICATION

In ReSum, an off-the-shelf LLM can serve as the summary tool. However, its role extends far beyond conventional conversation summarization. To guide web agents in persistent, goal-oriented exploration, the summary tool must perform logical reasoning over lengthy and noisy interaction histories, distill verifiable evidence from large text snippets, and propose actionable, well-scoped next steps grounded in web context. These capabilities typically exceed those of generic models that lack web-context reasoning, motivating the development of a specialized summary tool for ReSum.

To build an effective goal-oriented summary tool, we first conduct an empirical study comparing models of varying scales (Yang et al., 2024; Team, 2025b; OpenAI, 2025a). As shown in Figure 6 in Appendix E.1, smaller models often struggle to extract verifiable evidence from lengthy and noisy interaction histories, underscoring the importance of strong reasoning capabilities. While larger models excel in summarization, their high API costs and significant deployment overhead render them impractical. Therefore, we develop a smaller, deployable model that retains the goal-oriented summarization capabilities of larger models.

Development: We employ a rigorous quality control pipeline to develop our specialized summarization model. This process includes: (1) comprehensive teacher model evaluation and selection, (2) collection of 9K+ high-quality `<Conversation, Summary>` pairs from ReSum rollouts using the challenging **SailorFog-QA** benchmark (Li et al., 2025a), and (3) in-depth case analysis to

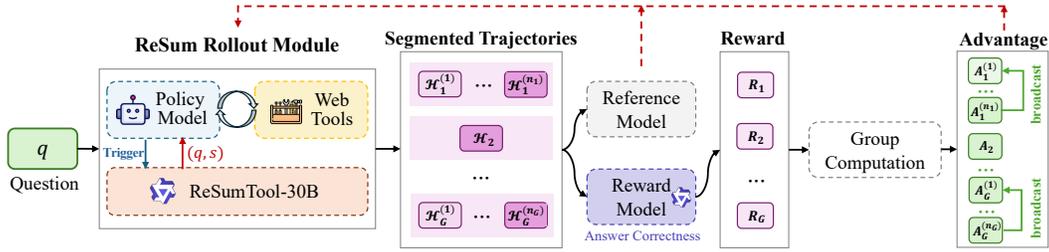


Figure 3: **Illustration of ReSum-GRPO.** ReSum periodically summarizes long trajectories and restarts from compressed states, resulting in segmented trajectories. A single trajectory-level reward is computed from the final answer, normalized within the group to obtain a trajectory-level advantage, and that advantage is **broadcast** to all segments within the same rollout.

ensure reliability. We distill this capability into Qwen3-30B-A3B-Thinking¹, selected for its MoE architecture that enables efficient deployment while maintaining strong reasoning capabilities. The resulting ReSumTool-30B demonstrates effective summarization performance, with detailed training configurations and empirical results provided in Appendix E.2.

3.3 RESUM-GRPO

The ReSum paradigm creates a novel query type $q' = (q, s)$ that combines the original user query q with a summary s . While agents can process such inputs, they may initially reason suboptimally from summarized contexts as this pattern was not encountered during training. Therefore, we employ RL to master such paradigm. Unlike supervised fine-tuning, which requires costly collection of expert-level ReSum trajectories and risks overwriting an agent’s existing skills, RL enables agents to adapt to this paradigm through self-evolution while retaining their inherent reasoning capabilities.

Trajectory Segmentation: The key modification of ReSum-RL is that ReSum naturally segments long trajectories into multiple episodes when summarization occurs. Consider a complete ReSum trajectory that undergoes K summarization events. This trajectory is naturally partitioned into the following $K + 1$ segments:

$$\begin{aligned} \mathcal{H}^{(1)} &= (q^{(0)}, \tau_1, a_1, o_1, \dots, \tau_{t_1}, a_{t_1}, o_{t_1}) \\ \mathcal{H}^{(2)} &= (q^{(1)}, \tau_{t_1+1}, a_{t_1+1}, o_{t_1+1}, \dots, \tau_{t_2}, a_{t_2}, o_{t_2}) \\ &\vdots \\ \mathcal{H}^{(K+1)} &= (q^{(K)}, \tau_{t_K+1}, a_{t_K+1}, o_{t_K+1}, \dots, \tau_T, a_T), \end{aligned}$$

where $q^{(0)} = q$ is the initial query, $q^{(k)} = (q, s^{(k)})$ is the compressed state after the k -th summarization, and a_T denotes the final answer. Each segment $\mathcal{H}^{(i)}$ forms an individual training episode with input $q^{(i-1)}$ and output $(\tau_{t_{i-1}+1}, a_{t_{i-1}+1}, o_{t_{i-1}+1}, \dots, \tau_{t_i}, a_{t_i})$. For trajectories that complete without summarization, we have the degenerate case $K = 0$, yielding a single segment that follows the same training format.

Reward Computation: To avoid manually designing per-segment rewards, we utilize a unified trajectory-level reward signal. From the last segment, we extract a_T and compute the reward as $R(a, a_T) \in \{0, 1\}$, where a represents the ground truth, using an LLM-as-Judge strategy (Gu et al., 2024; Li et al., 2024). This approach provides a single reward per complete trajectory, which can be shared across all its segments if necessary. Our reward design primarily focuses on answer correctness. Additionally, we perform **format checks** at each generation step: if the agent fails to adhere to specific tokens such as `<think>` `</think>`, the entire trajectory is terminated and assigned a zero reward. This binary format constraint ensures that the agent adheres to the required interaction pattern.

GRPO Integration (Figure 3): ReSum-RL modifies only the rollout collection by segmenting on summaries and adjusts the reward signal to be trajectory-level answer correctness. Consequently, it

¹<https://huggingface.co/Qwen/Qwen3-30B-A3B-Thinking-2507>

is compatible with various RL algorithms (Schulman et al., 2017; Christiano et al., 2017; Yu et al., 2025b). Specifically, we instantiate this with GRPO (Shao et al., 2024), resulting in ReSum-GRPO. For an initial question q , we sample a group of G rollouts, each producing n_g segments as $\{\mathcal{H}_g^{(i)}\}_{i=1}^{n_g}$. The objective can be written as:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{(q,a) \sim \mathcal{D}, \{\mathcal{H}_g^{(i)}\}_{g=1}^G, \{n_g\}_{g=1}^G \sim \pi_\theta} \left[\frac{1}{\sum_{g=1}^G n_g} \sum_{g=1}^G \sum_{i=1}^{n_g} \min \left(r_g^{(i)}(\theta) \hat{A}_g^{(i)}, \text{clip}(r_g^{(i)}(\theta), 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}}) \hat{A}_g^{(i)} \right) \right],$$

where $r_g^{(i)}(\theta)$ is the probability ratio for segment i in rollout g . In alignment with GRPO, rather than directly broadcasting rewards, we broadcast the trajectory-level advantage. For trajectory g , we extract the final answer $a_{g,T}$ from its last segment and compute a trajectory-level reward $R_g \in \{0, 1\}$. This reward is then normalized within the group to obtain the advantage as $\hat{A}_g = \frac{R_g - \text{mean}(\{R_1, \dots, R_G\})}{\text{std}(\{R_1, \dots, R_G\})}$, which is broadcast to all segments within rollout g as $\hat{A}_g^{(i)} = \hat{A}_g$ for $i \in \{1, \dots, n_g\}$. Such mechanism ensures a consistent learning signal per trajectory while leveraging GRPO’s group-wise stabilization.

In summary, the advantage broadcasting mechanism in ReSum-GRPO enables effective learning from all segments of long-horizon tasks by: (1) encouraging agents to reason successfully from compressed states, and (2) [ensuring early exploration steps receive appropriate bonus when they contribute to the final success](#). Notably, ReSum-GRPO only modifies long trajectories by utilizing segmented rollouts, while short trajectories are processed identically to standard GRPO. This design not only maintains training efficiency but also preserves the agent’s inherent reasoning patterns.

4 EXPERIMENTS AND ANALYSIS

4.1 EXPERIMENTAL SETUP

Benchmarks: To evaluate ReSum’s effectiveness in overcoming context limitations for complex queries, we conduct experiments on three challenging benchmarks where agents typically require extensive exploration: **GAIA** (Mialon et al., 2023), **BrowseComp-en** (Wei et al., 2025), and its Chinese counterpart **BrowseComp-zh** (Zhou et al., 2025a). For GAIA, we follow existing works by using the 103-sample text-only validation subset. We exclude simpler benchmarks such as SimpleQA (Wei et al., 2024), WebWalkerQA (Wu et al., 2025b), and xBench-DeepSearch (Xbench-Team, 2025), where most cases can be resolved within standard context limits, rendering the ReAct paradigm more suitable.

Evaluation: [Following standard practice in web agent research](#) (Gao et al., 2025; Dong et al., 2025), we consistently use Qwen2.5-72B-Instruct as the scoring model to assess whether the predicted answer aligns with the ground truth. Specifically, we report the average **Pass@1** over all test samples, as well as **Pass@3** across three rollouts for each sample. Unless otherwise stated, we set the maximum tool call budget to 60 for all inference paradigms to ensure a fair comparison.

Baselines & Implementations: We assess ReSum’s effectiveness under two settings: **training-free** and **training-required**. In the training-free setting, we directly apply the ReSum paradigm to various web agents without additional training. We compare its performance against ReAct and a simple baseline, **Recent History**, which maintains only the most recent $22k$ tokens of the conversation history by truncating older messages when context limit is exceeded. In the training-required setting, we compare ReSum-GRPO with the standard GRPO algorithm. Specifically, we evaluate both ReSum and ReAct paradigms on different RL-trained web agents to check whether ReSum-GRPO enhances the agents’ proficiency with ReSum. For ReSum inference, summarization is consistently triggered as the conversation history approaches the context limit, and leverages ReSumTool-30B unless specifically stated. Additionally, we include a comparison with **MEM1** (Zhou et al., 2025b), a representative context management method, under both settings, as detailed in [Section 4.4](#). Further implementation details are provided in Appendix D.

Table 1: **Performance comparison (in %) between paradigms under training-free settings. Bold** indicates results using ReSum with our developed ReSumTool-30B, which consistently outperforms ReAct and Recent History. **Blue highlights the best results for each backbone agent.** Results with [†] are sourced from Liu et al. (2025), representing leading pre-trained models paired with search and visit tools to illustrate the datasets’ performance landscape.

Agent	Paradigm	Summary Tool	GAIA		BrowseComp-zh		BrowseComp-en	
			Pass@1	Pass@3	Pass@1	Pass@3	Pass@1	Pass@3
Claude-4 [†]	ReAct	–	68.3	–	29.1	–	12.2	–
OpenAI-o3 [†]	ReAct	–	70.5	–	58.1	–	50.9	–
Kimi-K2 [†]	ReAct	–	57.7	–	28.8	–	14.1	–
DeepSeek-v3.1 [†]	ReAct	–	63.1	–	49.2	–	30.0	–
WebSailor-3B	ReAct	–	25.6	42.7	8.2	17.0	3.3	5.6
	Recent History	–	27.2	44.7	13.2	24.3	3.8	8.9
	ReSum	Qwen3-30B	27.5	45.6	6.9	14.5	4.2	7.8
		ReSumTool-30B	35.3	52.4	13.7	24.6	6.8	10.8
		GPT-OSS-120B	40.5	65.1	15.2	28.0	8.5	15.8
		Qwen3-235B	32.4	49.5	11.1	23.9	5.7	10.3
DeepSeek-R1-671B	39.2	60.2	13.0	23.5	7.5	13.4		
ReAct	–	31.7	44.7	13.2	25.6	5.7	10.3	
Recent History	–	33.0	48.5	15.2	28.0	5.2	9.4	
WebSailor-7B	ReAct	–	34.6	48.5	13.3	26.6	5.8	10.3
	ReSum	ReSumTool-30B	40.5	60.2	17.2	30.8	9.0	15.2
		GPT-OSS-120B	42.4	61.2	19.2	35.6	10.5	17.2
		Qwen3-235B	43.4	60.2	18.1	32.9	8.7	15.2
		DeepSeek-R1-671B	41.1	58.3	17.1	32.2	10.3	16.6
ReAct	–	45.0	60.2	23.9	38.4	12.8	21.8	
Recent History	–	40.1	56.3	24.1	40.1	10.3	16.7	
WebSailor-30B	ReAct	–	45.6	61.2	24.8	40.1	12.2	20.4
	ReSum	ReSumTool-30B	47.3	63.1	24.1	42.6	16.0	25.4
		GPT-OSS-120B	51.5	68.9	27.3	46.4	18.8	30.9
		Qwen3-235B	46.9	67.0	25.7	42.2	17.2	26.7
DeepSeek-R1-671B	49.2	71.8	27.1	41.5	13.7	22.6		

Choice of Web Agents: We conduct experiments on open-source web agents of varying scales to ensure a comprehensive evaluation, including **WebSailor-3B**², **WebSailor-7B**³, and **WebSailor-30B-A3B**⁴. Note that all these agents are constrained by a $32k$ token context limit.

4.2 PERFORMANCE OF THE TRAINING-FREE RESUM

Settings: We evaluate different inference paradigms on web agents, including **ReAct**, **ReSum**, and **Recent History**. All agents run under our unified inference framework with curated prompts. For summarization, we evaluate four off-the-shelf LLMs of varying scales, Qwen3-30B-A3B-Thinking (denoted as Qwen3-30B), GPT-OSS-120B, Qwen3-235B, and DeepSeek-R1-671B, alongside our developed ReSumTool-30B which leverages Qwen3-30B as the base. To contextualize performance, we also report results of leading pre-trained models like OpenAI-o3 (OpenAI, 2025b) and Kimi-K2 (Team et al., 2025) paired with search and visit tools. Quantitative results are presented in Table 1, revealing the following key findings⁵:

ReSum paradigm consistently outperforms ReAct due to extended exploration opportunities. The ReSum paradigm demonstrates statistically significant performance improvements across all agents and benchmarks, achieving average absolute gains of 4.5% over ReAct. This enhancement stems from ReSum’s ability to maintain coherent exploration through intelligent context compression, enabling agents to solve complex queries that would otherwise exceed context limits. While the

²<https://huggingface.co/Alibaba-NLP/WebSailor-3B>

³<https://huggingface.co/Alibaba-NLP/WebSailor-7B>

⁴This is a reproduced version using the same training data as the WebSailor series, with rejection sampling fine-tuning (RFT) applied to Qwen3-30B-A3B-Base model for 2 epochs.

⁵We attribute the performance gap between our reproduced WebSailor-3B/7B and the official WebSailor technical report (Li et al., 2025a) to implementation differences and prompt variations.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

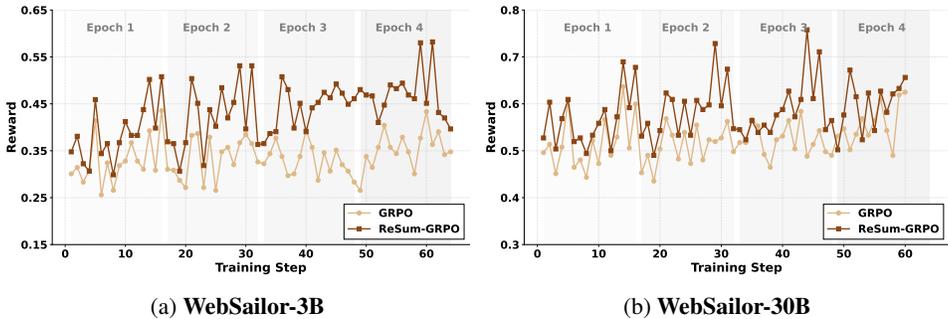


Figure 4: **Training dynamics comparison between GRPO and ReSum-GRPO.** ReSum-GRPO achieves higher rewards during training.

Recent History baseline also provides extended exploration, simple truncation disrupts context continuity and fails to preserve valuable information for continued reasoning.

For context summarization, our developed ReSumTool-30B achieves comparable performance to larger models while maintaining deployment efficiency. ReSumTool-30B consistently outperforms its base model Qwen3-30B across configurations when serving as the summary tool. Remarkably, ReSumTool-30B often matches or exceeds the performance of significantly larger models when used for summarization: on BrowseComp-zh with WebSailor-3B, it achieves 13.7% Pass@1, outperforming both Qwen3-235B (11.1%) and DeepSeek-R1-671B (13.0%) when they serve as summary tools. This demonstrates the effectiveness of our targeted training.

ReSum integration narrows the performance gap to several leading pre-trained models. Notably, WebSailor-30B with ReSumTool-30B realizes 16.0% Pass@1 on the BrowseComp-en benchmark, surpassing Claude-4-Sonnet (12.2%) and Kimi-K2 (14.1%). **While frontier models like OpenAI-o3 and DeepSeek-v3.1 maintain a performance lead, ReSum provides a practical path to enhance existing agents within computational constraints.**

ReSum paradigm further demonstrates **orthogonal effectiveness and broad compatibility** to stronger agents with extensive context windows, as supported by experiments on Tongyi-DeepResearch-30B-A3B (Team, 2025c) in **Appendix F**.

4.3 PERFORMANCE OF RESUM-GRPO

Settings: We use WebSailor models as the base for RL training, as they have already undergone RFT to acquire tool calling capabilities and provide a clean base without prior RL experience. For training data, we randomly select **1K samples** from the SailorFog-QA dataset (Li et al., 2025a) due to its high-quality and challenging characteristics. The training data scale is limited to 1K to reduce training costs while still enabling performance comparison. We emphasize that our goal is to **demonstrate the effectiveness of ReSum-GRPO** rather than pursuing performance limits through extensive training data. We compare ReSum-GRPO with standard GRPO, where trajectories are rolled out following the ReAct paradigm. Both RL algorithms are trained for 4 epochs with all hyper-parameters held consistent (details in Appendix D.2).

Training Dynamics: During each training step, the rewards on training data are shown in Figure 4, demonstrating that ReSum-GRPO consistently achieves higher rewards than GRPO. This is because ReSum mode extends the conversation of otherwise unsolvable questions for more exploration. As training progresses, ReSum-GRPO effectively encourages the agent to familiarize itself with this inference pattern, achieving higher rewards.

Overall Evaluation: We evaluate these RL-trained agents on both inference paradigms, along with their direct ReAct performance shown in Table 2. We also list the performance of powerful web agents trained with RL on **10K+** samples to contextualize the performance. We can conclude that: (1) **ReSum-GRPO successfully familiarizes agents with the ReSum paradigm, achieving more significant improvements across benchmarks.** For example, after ReSum-GRPO, WebSailor-3B achieves Pass@1 from 8.2% to 20.5% on BrowseComp-zh, demonstrating the effectiveness of RL training. (2) **The GRPO algorithm fails to enable agents to master summary-conditioned**

Table 2: **Performance comparison (in %) between RL algorithms.** ReSum-GRPO enables agents to become better familiarized with the ReSum paradigm, boosting performance. Results with [†] are sourced from Liu et al. (2025), representing powerful web agents trained with **10K+** samples.

Agent	RL	Paradigm	GAIA		BrowseComp-zh		BrowseComp-en	
			Pass@1	Pass@3	Pass@1	Pass@3	Pass@1	Pass@3
Qwen3-ARPO-14B	ARPO	ReAct	43.7	—	—	—	—	—
MiroThinker-8B [†] _{v0.1}	DPO	ReAct	46.6	—	13.6	—	8.7	—
MiroThinker-32B [†] _{v0.1}	DPO	ReAct	57.3	—	17.0	—	13.0	—
ASearcher-32B [†]	GRPO	ReAct	52.8	—	15.6	—	5.2	—
WebExplorer-8B [†]	GRPO	ReAct	50.0	—	32.0	—	15.7	—
WebSailor 3B	—	ReAct	25.6	42.7	8.2	17.0	3.3	5.6
	GRPO	ReAct	28.5	48.5	11.8	22.5	4.2	8.5
		ReSum	38.5	53.4	17.3	29.1	8.5	13.0
	ReSum-GRPO	ReSum	37.9	56.3	20.5	34.3	9.2	13.0
WebSailor 7B	—	ReAct	31.7	44.7	13.2	25.6	5.7	10.3
	GRPO	ReAct	34.0	47.6	18.7	31.8	5.8	10.0
		ReSum	37.2	53.4	25.4	40.8	8.5	15.0
		ReSum-GRPO	ReSum	42.4	60.2	27.1	39.5	12.3
WebSailor 30B	—	ReAct	45.0	60.2	23.9	38.4	12.8	21.8
	GRPO	ReAct	48.2	62.1	23.3	36.7	14.3	21.5
		ReSum	48.5	61.2	29.3	42.6	15.0	25.0
		ReSum-GRPO	ReSum	48.5	68.0	33.3	48.8	18.3

reasoning. GRPO is designed for familiarizing agents with the ReAct inference mode, which indeed boosts the agent’s ReAct inference, while applying the ReSum paradigm cannot outperform ReSum-GRPO trained counterparts by a significant margin, showing the necessity for paradigm adaptation. (3) **ReSum-GRPO enables agents to achieve performance comparable to agents trained with 10K+ samples.** Even compared with powerful open-source agents that have been trained for hundreds of steps on 10K+ samples, our ReSum-GRPO trained with only 1K+ samples enables the base agent to achieve comparable performance, e.g., WebSailor-30B achieves 33.3% on BrowseComp-zh, surpassing ASearcher-32B (15.6%) (Gao et al., 2025), MiroThinker-32B (17.0%) (Team, 2025a) and WebExplorer-8B (32.0%) (Liu et al., 2025).

Fine-grained Analysis (Appendix G): For **training efficiency**, ReSum-GRPO modifies only long trajectories by invoking summaries and restarting the conversation when the context window exceeds $32k$ tokens, which would otherwise lead to termination in GRPO. A comparison of training costs between the two algorithms is shown in Table 5 in Appendix, where ReSum-GRPO requires approximately 1.5x the training time of GRPO, an acceptable increase. For **inference efficiency**, we compare performance and resource consumption, including average token costs and tool call numbers, across different paradigms. As shown in Figures 8a and 8b in Appendix, ReSum paradigms achieve superior performance with reasonable resource utilization. For **qualitative analysis** (Appendix G.3), we present the full trajectories of three test cases on WebSailor-30B after ReSum-GRPO training. In one case, the agent directly derives the answer without using summaries, while in the other two, it successfully leverages summaries to generate the final answer. These examples demonstrate that ReSum-GRPO training not only familiarizes the agent with ReSum but also preserves its ability to directly derive answers within a few tool calls.

4.4 COMPARISON WITH MEM1

To situate ReSum within the landscape of context management approaches, we conduct the comparison with MEM1 (Zhou et al., 2025b), a representative method that addresses long-horizon reasoning through learned memory tokens with architectural modification.

Settings: We evaluate both paradigms under training-free and training-required settings using WebSailor-30B. For MEM1, we adapt its inference paradigm to web agents by replacing special tokens with `<think>` `</think>` and `<tool.call>` `</tool.call>` while preserving its core iterative memory consolidation mechanism. To ensure a fair comparison, both MEM1-GRPO and ReSum-GRPO utilize the GRPO algorithm (Shao et al., 2024) for optimization, with identical training

Table 3: Comparison between ReSum and MEM1 under both training-free and training-required settings using WebSailor-30B.

Setting	Paradigm	GAIA		BrowseComp-zh		BrowseComp-en	
		Pass@1	Pass@3	Pass@1	Pass@3	Pass@1	Pass@3
	ReAct	45.0	60.2	23.9	38.4	12.8	21.8
Training-free	MEM1	33.3	52.4	25.0	41.2	12.7	22.5
	ReSum	47.3	63.1	24.1	42.6	16.0	25.4
Training-required	MEM1-GRPO	35.7	54.4	29.1	45.0	19.5	29.7
	ReSum-GRPO	48.5	68.0	33.3	48.8	18.3	26.5

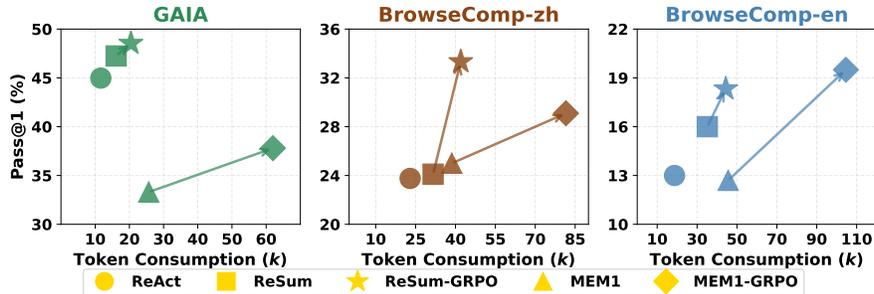


Figure 5: Average token consumption vs. performance across different paradigms. Token consumption refers to the total number of tokens in a complete trajectory for a query.

data and configurations. This ensures that the performance differences arise from paradigm design rather than the RL algorithm.

The results in Table 3 reveal several key findings: (1) **MEM1 exhibits weak compatibility with existing agents in the training-free setting.** Directly applying the MEM1 paradigm results in little to no performance improvement across datasets compared to ReAct, and in some cases, performance even declines. This is primarily due to MEM1’s inference paradigm deviating significantly from ReAct’s append-all-history approach, making it difficult for paradigm adaptation. (2) **Targeted RL training improves the performance.** For example, after MEM1-GRPO training, WebSailor-30B achieves a Pass@1 score of 19.5% on BrowseComp-en, demonstrating the benefits of specialized training for adapting to MEM1’s inference paradigm. (3) **MEM1 incurs substantial token costs due to its iterative reasoning process.** This paradigm integrates all thinking, planning, and memory into both inputs and outputs, leading to a substantial increase in token usage as shown in Figure 5. Specifically, MEM1-GRPO consumes nearly 3x more tokens than ReSum for a 1% improvement in Pass@1 score on BrowseComp-en. For a single query, a full trajectory can reach up to 110k tokens, resulting in enormous computational costs. In contrast, ReSum demonstrates better compatibility with existing agents and achieves a more favorable trade-off between cost and performance.

5 CONCLUSION

In this paper, we address the challenge that prohibits web agent to perform long-horizon searches: context limits. We propose ReSum, a novel inference paradigm that periodically employs summary tools for context summarization, enabling unbounded exploration. Furthermore, we introduce the tailored ReSum-GRPO algorithm, adapting agents to this inference paradigm through self-evolution. Extensive experiments on challenging benchmarks across agents demonstrate the effectiveness of both the ReSum paradigm and the targeted ReSum-GRPO training. Our **limitations** include reliance on external summary tools and rule-based summary invocation mechanisms. In future, we aim to equip agents with the capability to perform self-summarization and intelligently trigger summary calls at appropriate moments, eliminating dependence on external tools and predefined rules.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

ETHICS STATEMENT

This work adheres to the ICLR Code of Ethics. We confirm that no unauthorized datasets, test sets, or models were used in this research. All data utilized in this study were either publicly available or licensed for use. Additionally, we are committed to restricting the application of our model to avoid any harmful or unethical outcomes.

REPRODUCIBILITY STATEMENT

For all experimented web agents (WebSailor-3B to WebSailor-30B-A3B), we have provided links to their official release websites or included detailed reproduction information in the main text. For testing benchmarks, we strictly follow the official test set splits, with additional details outlined in the main text.

The complete training configurations for ReSumTool-30B are provided in Appendix E.2. Furthermore, all inference settings and the RL training details for both GRPO and ReSum-GRPO are thoroughly documented in Appendix D to ensure reproducibility. The curated prompts for both conversation summarization and summary-conditioned reasoning are detailed in Appendix C.

REFERENCES

- Anthropic. Introducing claude 4. *Technical report*, 2025. URL <https://www.anthropic.com/news/claude-4>.
- Hyunjoo Chae, Namyong Kim, Kai Tzu iunn Ong, Minju Gwak, Gwanwoo Song, Jihoon Kim, Sunghwan Kim, Dongha Lee, and Jinyoung Yeo. Web agents with world models: Learning and leveraging environment dynamics in web navigation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=moWiYJuSGF>.
- Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 4299–4307, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/d5e2c0adad503c91f91df240d0cd4e49-Abstract.html>.
- Guanting Dong, Hangyu Mao, Kai Ma, Licheng Bao, Yifei Chen, Zhongyuan Wang, Zhongxia Chen, Jiazhen Du, Huiyang Wang, Fuzheng Zhang, et al. Agentic reinforced policy optimization. *arXiv preprint arXiv:2507.19849*, 2025.
- Jiaxuan Gao, Wei Fu, Minyang Xie, Shusheng Xu, Chuyi He, Zhiyu Mei, Banghua Zhu, and Yi Wu. Beyond ten turns: Unlocking long-horizon agentic search with large-scale asynchronous rl, 2025. URL <https://arxiv.org/abs/2508.07976>.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L’elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. *ArXiv*, abs/2310.06825, 2023.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Serkan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. Search-rl: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025.
- Jina.ai. Jina, 2025. URL <https://jina.ai/>.
- Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. Llms-as-judges: A comprehensive survey on llm-based evaluation methods, 2024. URL <https://arxiv.org/abs/2412.05579>.
- Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, Weizhou Shen, Junkai Zhang, Dingchu Zhang, Xixi Wu, Yong Jiang, Ming Yan, Pengjun Xie, Fei Huang, and Jingren Zhou. Websailor: Navigating super-human reasoning for web agent. *arXiv preprint arXiv:2507.02592*, 2025a.
- Weizhen Li, Jianbo Lin, Zhuosong Jiang, Jingyi Cao, Xinpeng Liu, Jiayu Zhang, Zhenqiang Huang, Qianben Chen, Weichen Sun, Qiexiang Wang, Hongxuan Lu, Tianrui Qin, Chenghao Zhu, Yi Yao, Shuying Fan, Xiaowan Li, Tiannan Wang, Pai Liu, King Zhu, He Zhu, Dingfeng Shi, Piaohong Wang, Yeyi Guan, Xiangru Tang, Minghao Liu, Yuchen Eleanor Jiang, Jian Yang, Jiaheng Liu, Ge Zhang, and Wangchunshu Zhou. Chain-of-agents: End-to-end agent foundation models via multi-agent distillation and agentic rl, 2025b. URL <https://arxiv.org/abs/2508.13167>.

- 648 Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and
649 Zhicheng Dou. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint*
650 *arXiv:2501.05366*, 2025c.
- 651 Zhiyu Li, Shichao Song, Chenyang Xi, Hanyu Wang, Chen Tang, Simin Niu, Ding Chen, Jiawei
652 Yang, Chunyu Li, Qingchen Yu, Jihao Zhao, Yezhaohui Wang, Peng Liu, Zehao Lin, Pengyuan
653 Wang, Jiahao Huo, Tianyi Chen, Kai Chen, Kehang Li, Zhen Tao, Junpeng Ren, Huayi Lai,
654 Hao Wu, Bo Tang, Zhenren Wang, Zhaoxin Fan, Ningyu Zhang, Linfeng Zhang, Junchi Yan,
655 Mingchuan Yang, Tong Xu, Wei Xu, Huajun Chen, Haofeng Wang, Hongkang Yang, Wentao
656 Zhang, Zhi-Qin John Xu, Siheng Chen, and Feiyu Xiong. Memos: A memory os for ai system.
657 *arXiv preprint arXiv:2507.03724*, 2025d. URL <https://arxiv.org/abs/2507.03724>.
- 658 Junteng Liu, Yunji Li, Chi Zhang, Jingyang Li, Aili Chen, Ke Ji, Weiyu Cheng, Zijia Wu, Chengyu Du,
659 Qidi Xu, Jiayuan Song, Zhengmao Zhu, Wenhui Chen, Pengyu Zhao, and Junxian He. Webexplorer:
660 Explore and evolve for training long-horizon web agents, 2025. URL <https://arxiv.org/abs/2509.06501>.
- 661 Grégoire Mialon, Clémentine Fourier, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia:
662 a benchmark for general ai assistants. In *The Twelfth International Conference on Learning*
663 *Representations*, 2023.
- 664 OpenAI. gpt-oss-120b & gpt-oss-20b model card, 2025a. URL https://cdn.openai.com/pdf/419b6906-9da6-406c-a19d-1bb078ac7637/oai_gpt-oss_model_card.pdf.
- 665 OpenAI. Introducing openai o3 and o4-mini, 2025b. URL <https://openai.com/index/introducing-o3-and-o4-mini/>.
- 666 OpenAI. Openai deep research. <https://openai.com/index/introducing-deep-research/>, 2025.
- 667 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
668 optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- 669 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
670 Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical
671 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 672 Sijun Tan, Michael Luo, Colin Cai, Tarun Venkat, Kyle Montgomery, Aaron Hao, Tianhao Wu,
673 Arnav Balyan, Manan Roongta, Chenguang Wang, Li Erran Li, Raluca Ada Popa, and Ion Stoica.
674 rllm: A framework for post-training language agents. <https://pretty-radio-b75.notion.site/rLLM-A-Framework-for-Post-Training-Language-Agents-\21b81902c146819db63cd98a54ba5f31>, 2025. Notion Blog.
- 675 Zhengwei Tao, Jialong Wu, Wenbiao Yin, Junkai Zhang, Baixuan Li, Haiyang Shen, Kuan Li, Liwen
676 Zhang, Xinyu Wang, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. Webshaper: Agentically
677 data synthesizing via information-seeking formalization. *arXiv preprint arXiv:2507.15061*,
678 2025.
- 679 Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen,
680 Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong,
681 Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao,
682 Hongcheng Gao, Peizhong Gao, and et al. Kimi k2: Open agentic intelligence, 2025. URL
683 <https://arxiv.org/abs/2507.20534>.
- 684 MiroMind AI Team. Mirothinker: An open-source agentic model series trained for deep re-
685 search and complex, long-horizon problem solving. <https://github.com/MiroMindAI/MiroThinker>, 2025a.
- 686 Qwen Team. Qwen3 technical report, 2025b. URL <https://arxiv.org/abs/2505.09388>.
- 687 Tongyi DeepResearch Team. Tongyi-deepresearch. <https://github.com/Alibaba-NLP/DeepResearch>, 2025c.

- 702 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai
703 Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents.
704 *Frontiers of Computer Science*, 18(6):186345, 2024.
- 705
706 Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese,
707 John Schulman, and William Fedus. Measuring short-form factuality in large language models,
708 2024. URL <https://arxiv.org/abs/2411.04368>.
- 709 Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won
710 Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet
711 challenging benchmark for browsing agents. *arXiv preprint arXiv:2504.12516*, 2025.
- 712
713 Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu
714 Zhang, Zekun Xi, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. Webdancer: Towards
715 autonomous information seeking agency. *arXiv preprint arXiv:2505.22648*, 2025a.
- 716
717 Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Deyu Zhou,
718 Pengjun Xie, and Fei Huang. Webwalker: Benchmarking llms in web traversal. *arXiv preprint
arXiv:2501.07572*, 2025b.
- 719
720 Xbench-Team. Xbench-deepsearch, 2025. URL <https://xbench.org/agi/aisherech>.
- 721
722 Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe
723 Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents:
724 A survey. *Science China Information Sciences*, 68(2):121101, 2025.
- 725
726 Yuan Xia, Jingbo Zhou, Zhenhui Shi, Jun Chen, and Hai ting Huang. Improving retrieval augmented
727 language model with self-reasoning. In *AAAI Conference on Artificial Intelligence*, 2024.
- 728
729 Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. A-mem: Agentic
730 memory for llm agents. *arXiv preprint arXiv:2502.12110*, 2025.
- 731
732 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,
733 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint
arXiv:2412.15115*, 2024.
- 734
735 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
736 React: Synergizing reasoning and acting in language models. In *International Conference on
737 Learning Representations (ICLR)*, 2023.
- 738
739 Hongli Yu, Tinghong Chen, Jiangtao Feng, Jiangjie Chen, Weinan Dai, Qiyong Yu, Ya-Qin Zhang,
740 Wei-Ying Ma, Jingjing Liu, Mingxuan Wang, et al. Memagent: Reshaping long-context llm with
741 multi-conv rl-based memory agent. *arXiv preprint arXiv:2507.02259*, 2025a.
- 742
743 Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong
744 Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale.
745 *arXiv preprint arXiv:2503.14476*, 2025b.
- 746
747 W. Yu, Hongming Zhang, Xiaoman Pan, Kaixin Ma, Hongwei Wang, and Dong Yu. Chain-of-note:
748 Enhancing robustness in retrieval-augmented language models. In *Conference on Empirical
749 Methods in Natural Language Processing*, 2023.
- 750
751 Peilin Zhou, Bruce Leon, Xiang Ying, Can Zhang, Yifan Shao, Qichen Ye, Dading Chong, Zhiling
752 Jin, Chenxuan Xie, Meng Cao, et al. Browsecomp-zh: Benchmarking web browsing ability of
753 large language models in chinese. *arXiv preprint arXiv:2504.19314*, 2025a.
- 754
755 Zijian Zhou, Ao Qu, Zhaoxuan Wu, Sunghwan Kim, Alok Prakash, Daniela Rus, Jinhua Zhao, Bryan
756 Kian Hsiang Low, and Paul Pu Liang. Mem1: Learning to synergize memory and reasoning for
757 efficient long-horizon agents, 2025b. URL <https://arxiv.org/abs/2506.15841>.

A RELATED WORKS

Web Agents: Both proprietary and open-source communities have made significant strides in web agent development. Proprietary systems like DeepResearch (OpenAI, 2025) excel in complex web tasks but are hindered by closed architectures and inaccessible training data, limiting reproducibility and collaborative research. Open-source efforts, on the other hand, mainly focus on: data synthesis (e.g., data fuzzing in WebSailor and ASearcher), RL infrastructure, and algorithmic optimization (e.g., the specialized ARPO (Dong et al., 2025)). These advancements have propelled systems from addressing basic multi-hop question answering tasks to tackling more complex information-seeking challenges, such as the BrowseComp benchmark. Notable releases include WebSailor (Li et al., 2025a), WebShaper (Tao et al., 2025), ASearcher-QwQ-32B (Gao et al., 2025), and WebExplorer-8B (Liu et al., 2025). Despite these achievements, open-source agents remain fundamentally limited by the constrained exploration capabilities of the ReAct (Yao et al., 2023) paradigm, highlighting the need for new paradigms.

Context Management for Agents: The most widely used approach for context management in LLM-based agents is ReAct’s append-all-history strategy. While simple, this method leads to unbounded growth and rapid exhaustion, especially for complex queries. To address these issues, some methods introduce external components such as retrieval modules, e.g., A-Mem (Xu et al., 2025) and MemOS (Li et al., 2025d), to structure context more effectively. However, these solutions add significant computational overhead, increase system complexity, and integrate loosely with the agent. More recent approaches, such as MEM1 (Zhou et al., 2025b) and MemAgent (Yu et al., 2025a), allow agents to manage context internally through RL. **These methods innovate through architectural modification**, introducing learnable memory tokens that require end-to-end training of a new agent from scratch, which limits their applicability to pre-existing agents. In contrast, ReSum offers a lightweight **paradigmatic enhancement** to ReAct. This fundamental distinction yields key practical advantages including (1) **Training-free utility:** ReSum provides performance gains without any training, while methods like MEM1 can see performance drops in this setting; (2) **Data efficiency:** When training is desired, ReSum-GRPO achieves significant improvements using only 1K samples, reducing data cost; and (3) **Forward compatibility:** The decoupled summary tool can be independently improved, enhancing all ReSum-compatible agents without retraining.

Distinction from World-model-augmented Agents: While ReSum’s summarization shares the high-level goal of providing structured guidance with world models proposed in Chae et al. (2025), they address fundamentally different problems. World models focus on improving decision quality through dense, forward-looking planning at each step, whereas ReSum addresses context window exhaustion to enable decision continuation through sparse, backward-looking context compression. ReSum acts as an **occasional context manager** rather than an **integrated planner**, preserving compatibility with existing agents.

B ALGORITHM PSEUDO-CODE

In this section, we provide a detailed algorithmic description of the ReSum process in Algorithm 1.

Algorithm 1 ReSum Rollout with Periodic Context Summarization

Require: Query q , policy model π_θ , summary tool π_{sum} , tool environment \mathcal{R} , maximum tool calls B

Ensure: Final answer or failure

```

1: Initialize conversation history  $\mathcal{H}_0 \leftarrow (q)$ , tool call count  $b \leftarrow 0$ , round  $t \leftarrow 1$ 
2: while  $b < B$  do
3:   Generate reasoning and tool decision
4:    $(\tau_t, a_t) \sim \pi_\theta(\cdot \mid \mathcal{H}_{t-1})$   $\triangleright$  <think> </think> and <tool_call> </tool_call>
5:   if <answer> </answer> is detected in  $a_t$  then
6:     return final answer  $a_t$ 
7:   else if  $a_t$  is a tool call then
8:      $o_t \leftarrow \mathcal{R}(a_t)$   $\triangleright$  <tool_response> </tool_response>
9:      $\mathcal{H}_t \leftarrow \mathcal{H}_{t-1} \circ (\tau_t, a_t, o_t)$ 
10:     $b \leftarrow b + 1$ 
11:  else
12:    return failure  $\triangleright$  no answer or tool call
13:  end if
14:  if  $\text{Trig}(\mathcal{H}_t)$  then  $\triangleright$  Summarization trigger, e.g., token budget exceeded
15:     $s \sim \pi_{\text{sum}}(\cdot \mid \mathcal{H}_t)$   $\triangleright$  <summary> </summary> includes evidence and gaps
16:     $q' \leftarrow (q, s)$ ,  $\mathcal{H}_t \leftarrow (q')$   $\triangleright$  reset to compressed state
17:  end if
18:   $t \leftarrow t + 1$ 
19: end while
20: return failure  $\triangleright$  budget exhausted

```

C PROMPT

In this section, we provide the prompt used for invoking summary tools for context summarization within the ReSum paradigm. Note that we do not explicitly ask the summary tool to list current information gaps and provide clear action plans to avoid two potential issues: (1) the summary tool may lose focus on its primary task of consolidating key information, and (2) forced specification of information gaps may trap agents in cycles of repeated self-verification when answers are already within reach. However, we found that the summary tool **can intuitively and intelligently identify information gaps and suggest next-step plans when necessary**, demonstrating its emergent capability for strategic reasoning.

Prompt for Context Summarization

You are an expert at analyzing conversation history and extracting relevant information. Your task is to thoroughly evaluate the conversation history and current question to provide a comprehensive summary that will help answer the question.

Task Guidelines:

1. Information Analysis

- Carefully analyze the conversation history to identify truly useful information.
- Focus on information that directly contributes to answering the question.
- Do NOT make assumptions, guesses, or inferences beyond what is explicitly stated in the conversation.
- If information is missing or unclear, do NOT include it in your summary.

2. Summary Requirements

- Extract only the most relevant information that is explicitly present in the conversation.
- Synthesize information from multiple exchanges when relevant. Only include information that is certain and clearly stated in the conversation.
- Do NOT output or mention any information that is uncertain, insufficient, or cannot be confirmed from the conversation.

3. Output Format Your response should be structured as follows:

<summary>

- Essential Information: [Organize the relevant and certain information from the conversation history that helps address the question.]

</summary>

Strictly avoid fabricating, inferring, or exaggerating any information not present in the conversation. Only output information that is certain and explicitly stated.

Question {Question}

Conversation {Conversation History}

Please generate a comprehensive and useful summary.

After summary generation, we concatenate the initial question and the summary as a new formatted query for the agent to continue reasoning.

Prompt for Summary-conditioned Reasoning

Question {Question}

Below is a summary of the previous conversation. This summary condenses key information from earlier steps, so please consider it carefully. Assess whether the summary provides enough information to answer the question and use it as the basis for further reasoning and information gathering to answer the question.

Summary: {Summary}

D IMPLEMENTATION DETAILS

In this section, we elaborate on the implementation details of all inference paradigms and RL training configurations.

D.1 IMPLEMENTATION OF INFERENCE PARADIGMS

For our experimented agents, WebSailor-series, all are constrained by a context window of $32k$ tokens. We adopt the following settings for each inference paradigm. Note that for all inference paradigms, the maximum tool calling budget is 60 for a single query, and the LLM hyper-parameters are uniformly set with a `temperature` of 0.6 and `top_p` of 0.95.

- **ReAct:** Appending every thought, action, and observation into the conversation history. At each step, we detect the conversation window and terminate as failure if the agent has reached the context window without outputting the answer.
- **Recent History:** Whenever the context window has reached the limit, we truncate the conversation by only preserving the recent $22k$ tokens of messages, in this way, we can restart the conversation as well as preserving extra space for possible exploration.
- **MEM1:** Unlike ReAct’s append-all-history strategy, MEM1 maintains a constant context window, where the current query consists of the agent’s reasoning, planning, and the tool response from the previous turn. The agent then consolidates relevant information, generates a memory, and issues a tool call, iteratively refining the context to converge on the answer. For the **training-free** setting, we directly apply MEM1 inference to the web agent with prompt modifications. Specifically, to ensure compatibility with existing agents, we replace MEM1’s original special tokens, e.g., `<IS>`, `<query>`, with `<think>` `</think>` and `<tool_call>` `</tool_call>`. Additionally, the tool response from previous action is concatenated into the querying prompt, preserving the iterative structure of MEM1.
- **ReSum:** We consistently set the trigger for summarization as approaching the context limit, and then invoke ReSumTool-30B for conversation compression unless specifically stated.

In our current implementation of ReSum, we use a rule-based mechanism for summary triggering, i.e., approaching the token budget, which has the benefits of simple implementation and high efficiency as avoiding frequent summarization.

D.2 RL TRAINING CONFIGURATION

We implement both GRPO, ReSum-GRPO, and MEM1-GRPO for training web agents based on the `rLLM` framework (Tan et al., 2025). For these RL algorithms, all tool invocation results are excluded from loss calculation to prevent bias towards tool outputs following Jin et al. (2025); Dong et al. (2025).

Shared Hyper-parameters: For all RL algorithms, we consistently adopt a `batch_size` of 64, group size of 8, `learning_rate` of $2e - 6$, and 4 epochs due to the limited 1K training samples. Such consistent parameter settings ensure a fair comparison between algorithms.

Algorithm-specific Settings: For GRPO, the maximum number of tool calls is set to 40, with a total token limit of $32k$, where $2k$ tokens are allocated for the query prompt and $30k$ for responses, including thoughts, actions, and responses of tool calls. For ReSum-GRPO, the maximum number of tool calls is increased to 60, with $4k$ tokens allocated for the query prompt and $28k$ for responses. When the token limit is reached, the system invokes ReSumTool-30B to summarize the context, restart the conversation, and collect segmented trajectories from the prior process. For MEM1-GRPO, we adhere to MEM1 rollout process with trajectories optimized using the GRPO algorithm. Besides, we found MEM1’s paradigm difficult to apply to smaller models, as they frequently produced format errors and failed to follow complex memory consolidation instructions. This incompatibility disrupted RL training, preventing meaningful adaptation. Therefore, our MEM1 evaluation is limited to the stronger WebSailor-30B model.

E SUPPLEMENTARY MATERIALS FOR RESUMTOOL-30B

E.1 CASES OF LLMs IN CONTEXT SUMMARIZATION

Question: Among CS conferences, in 2025, which conference has exactly the same full paper submission deadline and the same CCF rank as IJCAI?

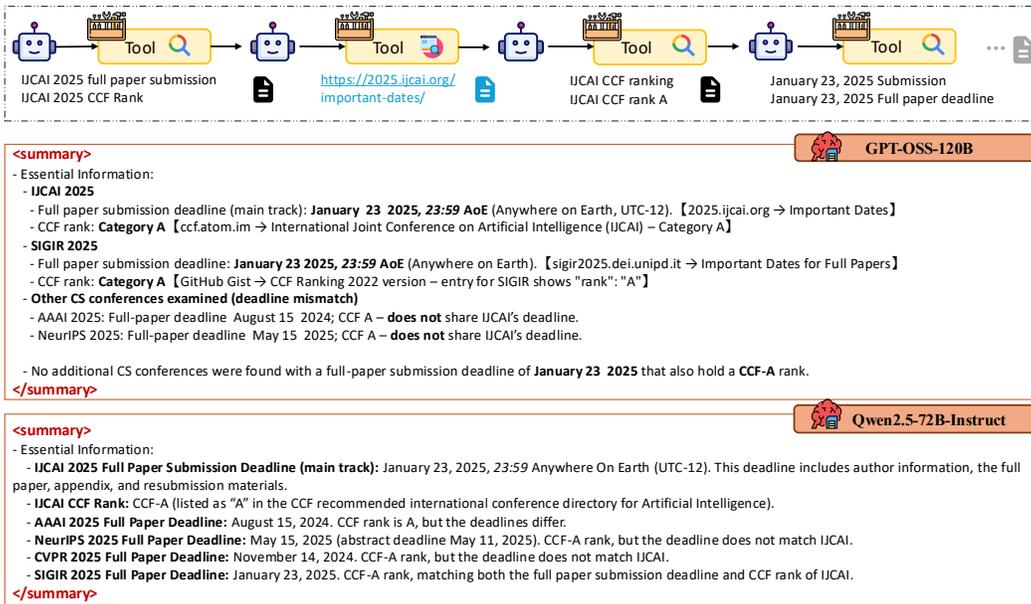


Figure 6: Comparison between summary contents generated by reasoning model GPT-OSS-120B (OpenAI, 2025a) and instruct model Qwen2.5-72B-Instruct (Yang et al., 2024).

We first conduct an empirical study comparing different models’ context summarization capabilities, including a reasoning model GPT-OSS-120B and an instruction model Qwen2.5-72B-Instruct.

Setting: The target question is “Among CS conferences, in 2025, which conference has exactly the same full paper submission deadline and the same CCF rank as IJCAI?”, with the ground truth answer being *SIGIR 2025*. We let a web agent perform ReAct inference on this case and truncate the conversation to the first 10 rounds of interaction, where the agent actively searches for related conferences and has gathered some information that can lead to the ground-truth answer. We then use the prompt in Appendix C to ask these two models to generate summaries, with output contents (highlighted parts aligned with the model’s original output in Markdown format) shown in Figure 6.

Observation: The comparison reveals significant differences in summarization quality and reasoning capabilities. GPT-OSS-120B demonstrates superior performance in several key aspects: (1) **structured organization** as it systematically categorizes information by conference with clear hierarchical formatting, (2) **comprehensive evidence gathering** as it identifies all relevant conferences and explicitly states why each candidate matches or fails the criteria, (3) **goal-oriented focus** as the summary directly addresses the question and highlights the final answer, and (4) **source attribution** as every piece of evidence is properly cited with specific sources. In contrast, Qwen2.5-72B-Instruct produces a more fragmented summary that lacks systematic organization. This highlights the necessity for specialized reasoning capabilities in context summarization tasks, especially in complex web search scenarios where structured evidence synthesis is essential for agent guidance.

E.2 TRAINING CONFIGURATIONS

In this subsection, we elaborate on the training process of ReSumTool-30B.

Data Collection: We collect $\langle \text{Conversation}, \text{Summary} \rangle$ pairs by performing ReSum rollout with WebSailor-30B as the agent and GPT-OSS-120B as the summary tool on a subset of the SailorFog-QA dataset. We select WebSailor-30B as the rollout model due to its zero API costs and

satisfactory search intelligence compared to other open-source LLMs. The summary tool is fixed to GPT-OSS-120B due to its high-quality summary generation and open-source availability. The dataset is selected for its difficulty, as SailorFog-QA mirrors challenging benchmarks like BrowseComp-en, where agents must utilize summary tools to solve problems. The collected summaries undergo format checking and are combined with the query prompts, including conversation history, to form over 9K $\langle \text{Input}, \text{Output} \rangle$ pairs. Here, Input represents the query prompt, while Output is the GPT-OSS-generated summary.

Training Hyper-parameters: We use Qwen3-30B-A3B-Thinking (Team, 2025b) as the base model and perform supervised fine-tuning on the collected data. The training configuration includes a `batch_size` of 64, 2 epochs, and a learning rate of 7×10^{-6} .

E.3 EVALUATION FOR RESUMTOOL-30B

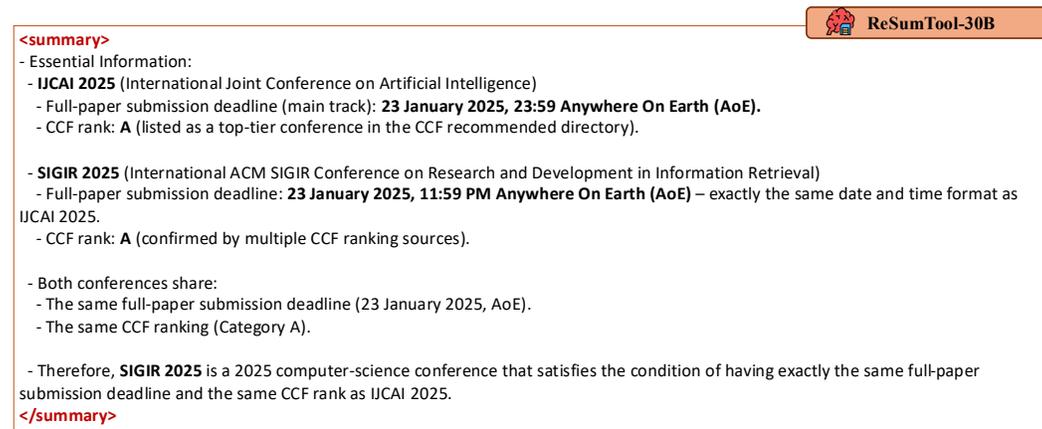


Figure 7: **Illustration of summary content generated by ReSumTool-30B with conversation and question mentioned in Figure 6.** The highlighted parts align with model’s original output in Markdown format.

To evaluate the performance of our trained model, we provide both quantitative results (Table 1) and qualitative analysis (Figure 7).

Quantitative Results: We measure the summary capability through agent ReSum inference performance, as agents rely on summaries to resume exploration. As analyzed in the main text, by comparing results with larger large reasoning models like DeepSeek-R1-671B and larger instruction models like Qwen3-235B, integrating our ReSumTool provides **comparable** performance boosts and significantly outperforms the Qwen3-30B Base, demonstrating its effectiveness.

Qualitative Analysis: We further provide summaries generated by ReSumTool-30B for illustration in Figure 7, where the solved question and conversation history exactly align with Figure 6. From this case, we can see that summaries generated by ReSumTool-30B exhibit reasonable structures, goal-focused organization, and comprehensive evidence gathering.

F APPLICABILITY TO STRONGER AGENTS WITH EXTENSIVE CONTEXT

Table 4: Performance comparison between ReAct and ReSum using Tongyi-DeepResearch-30B-A3B under training-free setting.

Context Limit	Tool Call	Paradigm	BrowseComp-zh		BrowseComp-en	
			Pass@1	Pass@3	Pass@1	Pass@3
32k	40	ReAct	41.2	57.4	27.7	43.1
		ReSum	43.8	62.3	34.5	53.3
48k	60	ReAct	42.5	58.5	32.8	48.7
		ReSum	46.7	62.6	38.2	54.5
64k	80	ReAct	43.6	60.9	36.3	52.4
		ReSum	48.6	66.1	40.3	57.8
96k	100	ReAct	46.0	62.3	39.8	56.5
		ReSum	47.9	66.4	41.0	57.0
128k	120	ReAct	45.7	62.3	42.2	59.2
		ReSum	46.6	62.6	44.5	59.5

In this section, we apply our ReSum paradigm to the powerful open-sourced web agent, **Tongyi-DeepResearch-30B-A3B** (Team, 2025c)⁶, which supports an extensive context size up to 128k tokens, to demonstrate **ReSum’s uniform compatibility and orthogonal effectiveness across agent capabilities**.

Settings: We conduct experiments across 5 context limit settings from 32k to 128k, with proportionally scaled tool call budgets (40 to 120 calls). For each setting, the ReAct paradigm terminates upon reaching the context limit, while ReSum invokes ReSumTool-30B for context summarization when approaching the limit. Notably, the experiments were conducted within our unified inference framework, which supports only the Search and Visit tools. Consistent with the model’s official inference parameters, we set the `temperature` to 0.85 and `top_p` to 0.95.

Results and Analysis: As demonstrated in Table 4, ReSum consistently outperforms ReAct across context sizes, confirming its effectiveness. The performance improvements are particularly substantial under smaller context constraints, highlighting the benefits of extended exploration enabled by summarization. Notably, even with extensive 128k context windows, ReSum still delivers meaningful gains, revealing that **complex problems demand exploration horizons beyond current context limits and that intelligent summarization remains crucial for tackling such challenges, aligning with established research** (Yu et al., 2023; Xia et al., 2024).

⁶<https://huggingface.co/Alibaba-NLP/Tongyi-DeepResearch-30B-A3B>

G SUPPLEMENTARY MATERIALS FOR EXPERIMENTS

In this section, we supplement the fine-grained experimental analysis of ReSum-GRPO, including training efficiency, inference costs, and concrete cases.

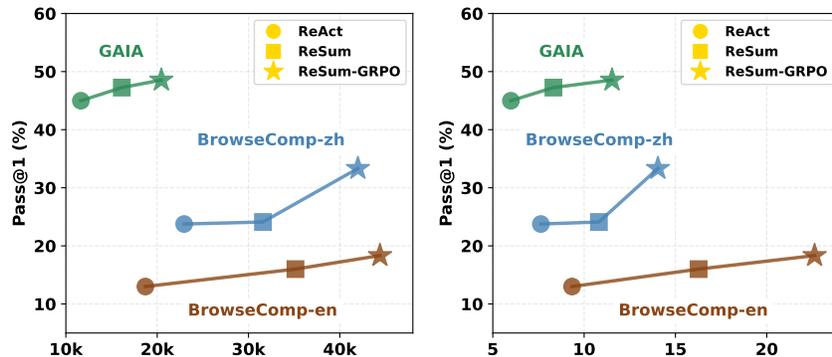
G.1 TRAINING EFFICIENCY

Table 5: Comparison of average time per single training step between RL algorithms.

Model	Device	GRPO	ReSum-GRPO
WebSailor-3B	8×144GB GPU _s	0.62 Hours	1.05 Hours
WebSailor-7B	8×144GB GPU _s	0.96 Hours	1.44 Hours
WebSailor-30B	16×144GB GPU _s	0.94 Hours	1.25 Hours

We provide the required devices and the average time for each training step for both RL algorithms in Table 5. Compared with GRPO, ReSum-GRPO modifies long trajectories by segmenting them based on summarization and then resumes the conversation for continued exploration. Consequently, the times required for both trajectory collection and policy model updates are lengthened. Based on the statistics in the table, ReSum-GRPO roughly increases training time by approximately 33% to 69% compared with GRPO, which is acceptable.

G.2 INFERENCE COSTS



(a) Number of cost tokens vs. Performance (b) Number of tool calls vs. Performance

Figure 8: **Resource consumption vs. performance across different paradigms.** We compare three paradigms: training-free ReAct, training-free ReSum, and ReSum-GRPO, consistently using WebSailor-30B. ReSum paradigms achieve higher performance with acceptable resource utilization across all benchmarks.

We further analyze resource consumption, i.e., the average number of tokens and tool calls required to correctly solve a query across different inference paradigms: training-free ReAct, training-free ReSum, and ReSum after ReSum-GRPO training. Token consumption refers to the total number of tokens in a complete trajectory for a query. Here, we only consider trajectories that successfully lead to a correct final answer. The results for WebSailor-30B across various benchmarks are presented in Figures 8a and 8b. From these results, we observe that in the training-free setting, ReSum significantly boosts performance with only marginal increases in resource costs compared to ReAct. Following targeted ReSum-GRPO training, agents become more inclined to rely on summaries for continued reasoning, which, while incurring additional resource costs, leads to even higher performance. Notably, ReSum paradigms achieve substantial performance improvements while maintaining resource costs within a reasonable range, e.g., typically $\sim 2x$ the original costs.

1188 G.3 CASE STUDY
1189

1190 In this subsection, we provide three representative cases of ReSum trajectories to demonstrate how
1191 this paradigm successfully helps solve questions. The agent is WebSailor-30B after ReSum-GRPO
1192 training. These three cases involve one that finishes without summarization, reflecting that ReSum-
1193 GRPO maintains the model’s capability to solve simpler questions in few tool calls, while the other
1194 two cases, from BrowseComp-en and BrowseComp-zh, respectively, show how the agent resumes
1195 reasoning based on summaries for successful question answering.

1196 **Case trajectory on BrowseComp-en: solved without summarization**
1197

1198 **Question:** As of 2023, tell me the scientific name of this carnivorous species based on the
1199 following clues: - its conservation status is ‘endangered’ - it is the only representative of its
1200 genus in the country - it is predominantly found in a region whose population density was
1201 132 per sqr km in 2011 - one of its local names roughly translates to “devil’s or ghost’s basket.”
1202

1203 **Round 1**1204 `<think>`

1205 The clues point to a very specific carnivorous species, but I still need to identify which
1206 one matches all the criteria: endangered status, being the only genus representative in its
1207 country, found in a region with a population density of 132 per sq km in 2011, and having a
1208 local name that means “devil’s or ghost’s basket.” By searching for combinations of these
1209 key phrases—especially focusing on the unique local name translation and the population
1210 density figure—I can cross-reference conservation lists, regional demographic data, and local
1211 terminology to pinpoint the species. Using multiple related queries increases the likelihood of
1212 surfacing authoritative sources or scientific papers that directly address this rare combination
1213 of traits. **I need to run this search to gather the necessary details for a definitive answer.**

1214 `</think>`1215 `<tool.call>`

1216 {"name": "search", "arguments": {"query": [{"carnivorous species endangered only
1217 representative genus country", "carnivorous species local names roughly translates devil’s
1218 ghost’s basket", "region population density 132 per sqr km 2011", "endangered carnivorous
1219 species region population density 132 per sqr km 2011"}]}}

1220 `</tool.call>`1221 `<tool.response>`

1222 A Google search for “carnivorous species endangered only representative genus country”
1223 found 10 results:

1224 1. [Conservation of carnivorous plants in the age of extinc-
1225 tion](https://www.sciencedirect.com/science/article/pii/S2351989420308131)

1226 CR species were recorded from 17 countries, although Brazil, Indonesia and the Philippines
1227 together harboured over half of all CR species (Table 2). Table 1.

1228 2. [List of countries harbouring highest numbers of threatened
1229 ...](https://www.researchgate.net/figure/List-of-countries-harboursing-highest-numbers-of-
1230 threatened-carnivorous-plant-species_tbl2_344197206)

1231 Around 860 unique species of flesh-eating plants are reported worldwide out of which 5
1232 genera belong to 3 families of carnivorous plants are recorded in India.

1233 ...

1234 `</tool.response>`1235 **Round 2**1236 `<think>`

1237 The connection between a carnivorous species with a local name meaning “devil’s or ghost’s
1238 basket” and its status **as the only genus representative in India** is becoming clearer, espe-
1239 cially since Meghalaya—a region with a population density of 132 per sq km in 2011—has
1240 been identified as a possible location. However, I still need to **confirm which carnivorous
1241 plant in India fits all these criteria**, particularly focusing on its scientific name and whether
it is indeed the sole genus representative. By searching for terms like “carnivorous plant India

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

only genus representative” and “carnivorous plant India local names devil’s basket”, I can directly target authoritative sources or botanical references that mention both the unique local name and the taxonomic status. I need to run this search now to pinpoint the exact scientific name and verify all the clues.

</think>

<tool_call>

```
{“name”: “search”, “arguments”: {“query”: [“carnivorous plant India only genus representative”, “carnivorous plant India local names devil’s basket”]}}
```

</tool_call>

<tool_response>

A Google search for “carnivorous plant India only genus representative” found 10 results:

- [First report of the chloroplast and mitochondrial genomes of the ...](https://bmresnotes.biomedcentral.com/articles/10.1186/s13104-025-07082-w)
Nepenthes, sometimes known as tropical pitcher plants or monkey cups, is a carnivorous plant genus that contains more than 160 species. *Nepenthes khasiana*, India’s sole representative of the genus, is a rare and endangered dioecious plant endemic to North-east India.
- [[PDF] *Nepenthes khasiana* - IUCN Red List](https://www.iucnredlist.org/species/pdf/4900..) Date published: Jul 16, 2014
Nepenthes khasiana, the only representative of the genus *Nepenthes* in India, belongs to monotypic family Nepenthaceae (Mandal and Mukherjee 2011)...

...

</tool_response>

Round 3

<think>

The Wikipedia page for *Nepenthes khasiana* is likely to provide authoritative and consolidated information on its conservation status, taxonomic uniqueness in India, and native distribution—all of which are crucial for confirming that this species matches every clue in the question. Additionally, Wikipedia often includes sections on local names and their meanings, which will let me verify whether “Memang-koksi” or similar terms translate to “devil’s basket” or “ghost’s basket.” By visiting this page, I can efficiently extract all the required details in one place and ensure that the species fits every aspect of the query, including its association with Meghalaya, a region with the specified population density. I need to **review this page now to confirm and synthesize these key facts.**

</think>

<tool_call>

```
{“name”: “visit”, “arguments”: {“url”: [“https://en.wikipedia.org/wiki/Nepenthes_khasiana”], “goal”: “Confirm conservation status, genus representation, distribution region, and local names translating to devil’s basket or ghost’s basket”}}
```

</tool_call>

<tool_response>

The useful information in https://en.wikipedia.org/wiki/Nepenthes_khasiana for user goal Confirm conservation status, genus representation, distribution region, and local names translating to devil’s basket or ghost’s basket as follows: ...

</tool_response>

Round 4

<think>

Although I have strong evidence that Meghalaya is the region with a population density of 132 per sq km in 2011 and that *Nepenthes khasiana* is the only *Nepenthes* species native to India, I want to ensure **there is no ambiguity by directly confirming Meghalaya’s 2011 population density** figure from an authoritative source. The Wikipedia page listing Indian states and union territories by population should provide official census data, allowing me to verify this crucial demographic detail and solidify the link between the species’ distribution and the region in question. I need to visit this page now to extract the exact 2011 population

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

density for Meghalaya.

</think>

<tool_call>

{“name”: “visit”, “arguments”: {“url”: [“https://en.wikipedia.org/wiki/List_of_states_and_union_territories_of_India_by_population”], “goal”: “Confirm population density of Meghalaya in 2011”}}

</tool_call>

<tool_response>

The useful information in https://en.wikipedia.org/wiki/List_of_states_and_union_territories_of_India_by_population for user goal Confirm population density of Meghalaya in 2011 as follows: ...

</tool_response>

Round 5

<think>

I determined that the scientific name is *Nepenthes khasiana* by systematically matching each clue in the question to authoritative sources. First, I **confirmed** that Meghalaya is the only Indian region with a 2011 population density of 132 per sq km, directly from census data. Then, I **verified** that *Nepenthes khasiana* is an endangered carnivorous pitcher plant, the sole representative of its genus in India, and is endemic to the Khasi, Jaintia, and Garo Hills—regions within Meghalaya. Finally, I **found** multiple reputable references stating that the Garo people call it “memang-koksi”, which literally means “basket of the devil”, aligning with the local name clue. By cross-referencing these facts, it became clear that *Nepenthes khasiana* is the only species fitting all the criteria, so I am confident in this answer.

</think>

<answer>

The scientific name of the carnivorous species described is *Nepenthes khasiana*.

</answer> ✓

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Case trajectory on BrowseComp-zh (translated into English): resumed from summary

Question: A renowned film from the last century, adapted from a film literary script whose title (different from the film title) contains numbers. After more than 60 years, it was restored and re-screened. The film’s director worked as a teacher with a four-character female actress in southwestern China during the 1930s-40s. What is the birth place of this female actress?

Round 1

`<think>`

To identify this film and its literary script, I need to find a work from the last century that was adapted from a film literary script containing numbers, and was restored and re-screened after more than 60 years. Since the question mentions “script title differs from film title” and “contains numbers”, this provides unique clues for retrieval. By searching for “last century renowned film adapted film literary script numbers”, “film literary script numbers 60+ years restoration re-screening”, “renowned film adapted script numbers 60+ years”, and “film adapted film literary script 60+ years restoration”, I can maximize coverage of relevant works and events, helping to quickly identify the qualifying film and its script. **Next, I need to search these keywords to obtain specific film and script names and related background.**

`</think>`

`<tool_call>`

```
{ "name": "search", "arguments": { "query": ["last century renowned film adapted film literary script numbers", "film literary script numbers 60+ years restoration re-screening", "renowned film adapted script numbers 60+ years", "film adapted film literary script 60+ years restoration"] } }
```

`</tool_call>`

`<tool_response>`

A Google search for “last century renowned film adapted film literary script numbers” found 10 results: 1. ...

`</tool_response>`

Round 2

`<think>`

Currently, I need to **further identify the specific information about that renowned film from the last century**, especially the unique clue that “script title differs from film title” and “contains numbers”. The previous search involved keywords like “numbers” and “restoration”, but the results are still scattered and haven’t directly pointed to a film meeting all conditions. Therefore, I decide to refine the search further, focusing on “last century film adapted script numbers”, “film literary script numbers 60+ years restoration”, “renowned film literary script numbers restoration” and other keywords, which can more precisely filter out film literary scripts related to numbers and pay attention to cases of restoration and re-screening after 60+ years. This can directly narrow the scope and find the most likely target film and its director and actress clues.

`</think>`

`<tool_call>`

```
{ "name": "search", "arguments": { "query": ["last century film adapted script numbers", "film literary script numbers 60+ years restoration", "renowned film literary script numbers restoration"] } }
```

`</tool_call>`

`<tool_response>`

A Google search for “last century film adapted script numbers” found 10 results ...

`</tool_response>`

Round 3

`<think>`

“Shangganling” film has appeared multiple times in previous searches, and it’s clearly mentioned that it was adapted from the film literary script “**Twenty-Four Days**”, with the

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

script title containing numbers, and the film had a 4K restoration in 2020. Now **I need to confirm who directed “Shangganling” and whether the director worked as a teacher with a four-character female actress in southwestern China during the 1930s-40s**, to identify this actress and her birth place. For this purpose, I decide to search with keywords like “Shangganling director”, “Shangganling director 30 40 teacher”, “Shangganling director southwest teacher”, “Shangganling director four-character female actress” to precisely obtain the director’s historical connection with the actress and their detailed background.

</think>

<tool_call>

```
{ "name": "search", "arguments": { "query": ["Shangganling director", "Shangganling director 30 40 teacher", "Shangganling director southwest teacher", "Shangganling director four-character female actress"] } }
```

</tool_call>

<tool_response>

A Google search for “Shangganling director” found 10 results ...

</tool_response>

Round 4

<think>

Currently, I have grasped the connection between “Shangganling” film and **director Sha Meng**, but there’s still a lack of direct evidence about whether Sha Meng worked as a teacher with a four-character female actress in southwestern China during the 1930s-40s. To identify this actress, the most effective way is to focus on keywords like “Sha Meng teacher southwest”, “Sha Meng teacher 30 40”, “Sha Meng southwest teacher female actress”, “Sha Meng teacher female actress” for searching, which can precisely screen historical records or biographies of Sha Meng working as a teacher with actresses in the southwest region during that period. If relevant materials can be found, we can further trace the actress’s birth place.

</think>

<tool_call>

```
{ "name": "search", "arguments": { "query": ["Sha Meng teacher southwest", "Sha Meng teacher 30 40", "Sha Meng southwest teacher female actress", "Sha Meng teacher female actress"] } }
```

</tool_call>

<tool_response>

A Google search for “Sha Meng teacher southwest” found 10 results ...

</tool_response>

Round 5

<think>

Currently, about whether Sha Meng worked as a teacher with a four-character female actress in southwestern China during the 1930s-40s, there’s still a lack of direct evidence. To identify this actress, the most effective way is to focus on keywords like “Sha Meng teacher southwest”, “Sha Meng teacher 30 40”, “Sha Meng southwest teacher female actress”, “Sha Meng teacher female actress” for searching, which can precisely screen historical records or biographies of Sha Meng working as a teacher with actresses in the southwest region during that period. **If relevant materials can be found, we can further trace the actress’s birth place.**

</think>

<tool_call>

```
{ "name": "search", "arguments": { "query": ["Sha Meng teacher southwest", "Sha Meng teacher 30 40", "Sha Meng southwest teacher female actress", "Sha Meng teacher female actress"] } }
```

</tool_call>

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

<tool_response>

A Google search for “Sha Meng teacher southwest” found 10 results ...

</tool_response>

Round 6

<think>

Currently, about whether Sha Meng worked as a teacher with a four-character female actress in southwestern China during the 1930s-40s, the direct evidence is still not clear enough, but I notice that the name “Ouyang Ruqiu” appears multiple times in related historical materials and has close connections with Sha Meng. Considering the limitation of “four-character female actress”, combined with the history of Sha Meng and Ouyang Ruqiu working as art teachers in the southwest region during the Anti-Japanese War, this is very likely the actress referred to in the question. **To confirm whether Ouyang Ruqiu was indeed Sha Meng’s collaborator when he worked as a teacher in the southwest, and to further obtain her birth place and other detailed information, I need to search with keywords** like “Sha Meng Ouyang Ruqiu teacher”, “Ouyang Ruqiu teacher southwest”, “Sha Meng Ouyang Ruqiu”, “Ouyang Ruqiu birth place” for precise retrieval.

</think>

<tool_call>

```
{ "name": "search", "arguments": { "query": ["Sha Meng Ouyang Ruqiu teacher", "Ouyang Ruqiu teacher southwest", "Sha Meng Ouyang Ruqiu", "Ouyang Ruqiu birth place"] } }
```

</tool_call>

<tool_response>

A Google search for “Sha Meng Ouyang Ruqiu teacher” found 10 results ...

</tool_response>

Round 7

Question: ...

<summary>

- The film “Shangganling” (produced by Changchun Film Studio in 1956) was adapted from the film literary script “Twenty-Four Days”. **The script title differs from the film title and contains numbers.** The film underwent 4K digital restoration and re-screening around 2020, more than 60 years after its original release.
- One of the directors of “Shangganling”, Sha Meng (original name Liu Shangwen, 1907-11-02 to 1964-06-26), worked as a teacher alongside actress Ouyang Ruqiu at the Yucai School founded by Tao Xingzhi in Beibei, Chongqing during the Anti-Japanese War period (1939). The two later married.
- Ouyang Ruqiu (born 1918 to December 27, 2007) was **a four-character name actress**, performing artist, and educator, **born in Xiaoxian County, Anhui Province.** Therefore, the birthplace of the four-character name actress referred to in the question is Xiaoxian County, Anhui Province.

</summary>

<think>

The summary has already provided complete clues and a clear answer, including the film “Shangganling”, the connection between its director Sha Meng and actress Ouyang Ruqiu, as well as Ouyang Ruqiu’s birthplace. **Since all key information has been verified and the logical chain is complete, I do not need to conduct additional searches or verification. Therefore, I can directly provide the final answer.**

</think>

<answer>

Based on the detailed clues and multiple rounds of verification, the birth place of the female actress Ouyang Ruqiu involved in the question is Xiaoxian County, Anhui Province.

</answer> ✓

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

Case trajectory on BrowseComp-en: resumed from summary

Question: There exists a peculiar study comparing certain abilities of two distinct species (species A and species B) done in the years between 2000 and 2005, inclusive. For the research experiment, 450 specimens of each species were used. There were three researchers from the university founded between 1800 and 1825, inclusive, who worked on this paper. In the years between 2005 and 2010, inclusive, this research earned recognition for its ability to amuse before inspiring deeper thought. Within the results and discussion of the study comparing certain abilities of two distinct species, what were the recorded measurements in cm for 50% of Species A and Species B, respectively, after the linearisation of the curves? Below is a summary of the previous conversation. This summary condenses key information from earlier steps, so please consider it carefully. Assess whether the summary provides enough information to answer the question and use it as the basis for further reasoning and information gathering to answer the question.

<summary>

- Essential Information: The study comparing jump performances of two flea species (*Ctenocephalides canis* and *Ctenocephalides felis felis*) used 450 specimens per species. After linearisation of the cumulative jump-height curves, the mean height jump performed by 50% of the fleas was recorded as **15.5cm** for *Ctenocephalides felis felis*. **The corresponding value for *Ctenocephalides canis* is not fully provided in the available excerpt.** The research falls within the 2000-2005 time frame and received recognition for its amusing yet thought-provoking nature between 2005 and 2010.

</summary>

Round 1

<think>

The summary confirms the study is about flea jump performance, with 450 specimens per species and the mean jump height for 50% of *C. felis felis* after linearisation as 15.5 cm, but the corresponding value for *C. canis* is missing. To fully answer the question, **I need to find the exact measurement for 50% of *C. canis* after linearisation, as well as verify the details about the researchers and the university, and confirm the recognition period.** By searching for combinations of terms like “*Ctenocephalides canis* 50% linearisation cm” and including references to the number of specimens, university founding date, and recognition period, I can efficiently target sources that mention both species’ results, the study’s authorship, and its subsequent recognition. This approach should surface either the original study or authoritative summaries that provide the missing data point and corroborate the context. I need to run this search now to fill in the remaining gaps and ensure a complete, accurate answer.

</think>

<tool.call>

```
{“name”: “search”, “arguments”: {“query”: [“Ctenocephalides canis 50% linearisation cm”, “Ctenocephalides canis Ctenocephalides felis felis 450 specimens”, “Ctenocephalides canis university founded between 1800 and 1825”, “Ctenocephalides canis recognition amuse inspiring deeper thought 2005 2010”]}}
```

</tool.call>

<tool.response>

A Google search for “*Ctenocephalides canis* 50% linearisation cm” found 10 results:

1. [A comparison of jump performances of the dog flea ... - PubMed](https://pubmed.ncbi.nlm.nih.gov/10962162/) The mean height jump carried out by 50% of fleas was calculated after linearisation of the curves: it was 15.5 and 13.2cm for *C. canis* and *C. felis*, ...
2. [*Ctenocephalides (felis) felis* (cat flea) *Ctenocephalides (felis) canis* ...](https://wcv.m.usask.ca/learnaboutparasites/parasites/ctenocephalides.php) Flea larvae are highly mobile and are able to move up to 40 cm on a suitable substrate. They feed on blood-rich fecal material from the adult fleas, on non- ...

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

`</tool_response>`

Round 2

`<think>`

The PubMed link appears to lead directly to the original study comparing the jump performances of *C. canis* and *C. felis felis*, and the search snippets already mention both species' mean heights after linearisation of the curves. However, to ensure accuracy and completeness—especially regarding the exact values for 50% of each species, the number of specimens used, the university affiliation of the researchers, and any mention of recognition or awards—I need to examine the full abstract or main text of the article. By visiting this page, I can extract the authoritative, primary-source data required to answer the question precisely and confirm all contextual details. **I need to visit this PubMed page now to gather these specifics.**

`</think>`

`<tool_call>`

```
{“name”: “visit”, “arguments”: {“url”: [“https://pubmed.ncbi.nlm.nih.gov/10962162/”],
“goal”: “Extract the mean height jump carried out by 50% of fleas after linearisation of the
curves for C. canis and C. felis felis, number of specimens per species, and details of the
university affiliation of the three researchers” }}
```

`</tool_call>`

`<tool_response>`

The useful information in <https://pubmed.ncbi.nlm.nih.gov/10962162/> for user goal Extract the mean height jump carried out by 50% of fleas after linearisation of the curves for *C. canis* and *C. felis felis*, number of specimens per species, and details of the university affiliation of the three researchers as follows: ...

`</tool_response>`

Round 3

`<think>`

The details about the university founding date and the recognition period are still unclear, and confirming these is crucial to fully satisfy the question's constraints. The previous evidence established the species and the measurements, but not the institutional affiliation or the timing of the study's recognition. By searching for combinations like “*Ctenocephalides canis* university founded between 1800 and 1825” and “*Ctenocephalides canis* recognition amuse inspiring deeper thought 2005 2010”, I can directly target sources that mention both the university's founding era and any awards or notable recognition within the specified timeframe. This approach should surface either the original publication, institutional profiles, or news articles that tie the study to a qualifying university and document its subsequent recognition. **I need to run this search now to bridge these remaining gaps and ensure a complete, authoritative answer.**

`</think>`

`<tool_call>`

```
{“name”: “search”, “arguments”: {“query”: [“Ctenocephalides canis university founded
between 1800 and 1825”, “Ctenocephalides canis recognition amuse inspiring deeper thought
2005 2010”]}}
```

`</tool_call>`

`<tool_response>`

A Google search for “*Ctenocephalides canis* university founded between 1800 and 1825” found 10 results:

1. [EENY589/IN1033: Dog Flea *Ctenocephalides canis* (Curtis) (Insecta)](<https://edis.ifas.ufl.edu/publication/IN1033>) Missing: 1800 1825
2. [*Ctenocephalides canis* is the dominant flea species of dogs in the ...](<https://pmc.ncbi.nlm.nih.gov/articles/PMC5859511/>)...

`</tool_response>`

Round 4

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

`<think>`

The exact number of specimens used in the study is a critical detail for fully validating the match to the question’s constraints, especially since the abstract and PubMed summary referenced group sizes of 10 fleas but did not clarify whether this translated to 450 specimens per species overall. The ScienceDirect article is likely to contain the full methods section or a table with the precise number of specimens used in the experiment, which is essential for confirming this key point. By visiting this page, I can directly verify the sample size and ensure the study aligns with the requirement of 450 specimens per species. **I need to examine this article now to extract the definitive specimen count and finalize the answer.**

`</think>`

`<tool_call>`

```
{“name”: “visit”, “arguments”: {“url”: [“https://www.sciencedirect.com/science/article/abs/pii/S0304401700002740”], “goal”: “Extract number of specimens used in the study comparing jump performances of C. canis and C. felis felis”}}
```

`</tool_call>`

`<tool_response>`

The useful information in <https://www.sciencedirect.com/science/article/abs/pii/S0304401700002740> for user goal Extract number of specimens used in the study comparing jump performances of C. canis and C. felis felis as follows: ...

`</tool_response>`

Round 5

`<think>`

The precise institutional affiliation of the researchers is a crucial detail for fully satisfying the question’s requirements, especially since the university’s founding date between 1800 and 1825 is a key constraint. While I have already identified the study’s content and measurements, I still lack confirmation of the university involved and whether it matches the founding date criteria. By searching for “Ctenocephalides canis study university”, I am targeting sources that are likely to mention the authors’ institutional affiliations, either in the abstract, full text, or in related summaries. This should help me directly identify the university and verify if it fits the required timeframe, thereby closing the remaining gap in the answer. **I need to run this search now to obtain the definitive institutional information.**

`</think>`

`<tool_call>`

```
{“name”: “search”, “arguments”: {“query”: [“Ctenocephalides canis study university”]}}
```

`</tool_call>`

`<tool_response>`

A Google search for “Ctenocephalides canis study university” found 10 results:

1. [Ctenocephalides canis is the dominant flea species of dogs in the ...](<https://pmc.ncbi.nlm.nih.gov/articles/PMC5859511/>) Date published: Mar 20, 2018 In the present study, we investigated the flea infestation among 116 outdoor dogs (57 females and 59 males) in 8 rural areas of Jeonnam Province ...
2. [EENY589/IN1033: Dog Flea Ctenocephalides canis (Curtis) (Insecta)](<https://edis.ifas.ufl.edu/publication/IN1033>) This document provides an overview of the dog flea, Ctenocephalides canis, including its biology, distribution, and medical significance....

`</tool_response>`

Round 6

`<think>`

The details about the university affiliation and the recognition period remain ambiguous, and these are critical to fully satisfy the question’s constraints. The previous evidence confirms the species, measurements, and sample sizes, but not the institutional connection or the timing of any notable recognition. By searching for combinations like “Ctenocephalides canis university founded between 1800 and 1825”, “Ctenocephalides canis recognition amuse

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

inspiring deeper thought 2005 2010”, and “Ctenocephalides canis 450 specimens university”, I am directly targeting sources that are likely to mention both the university’s founding era and any awards or recognition within the specified timeframe. This approach should surface either the original publication, institutional profiles, or news articles that tie the study to a qualifying university and document its subsequent recognition. I need to run this search now to bridge these remaining gaps and ensure a complete, authoritative answer.

</think>

<tool_call>

```
{“name”: “search”, “arguments”: {“query”: [“Ctenocephalides canis university founded between 1800 and 1825”, “Ctenocephalides canis recognition amuse inspiring deeper thought 2005 2010”, “Ctenocephalides canis 450 specimens university”]}}
```

</tool_call>

<tool_response>

A Google search for “Ctenocephalides canis university founded between 1800 and 1825” found 10 results:

1. [EENY589/IN1033: Dog Flea Ctenocephalides canis (Curtis) (Insecta)(https://edis.ifas.ufl.edu/publication/IN1033) Missing: 1800 1825
2. [Ctenocephalides canis is the dominant flea species of dogs in the ...](https://pmc.ncbi.nlm.nih.gov/articles/PMC5859511/) Missing: 1800 1825

...

</tool_response>

Round 7

<think>

The exact institutional affiliation of the researchers is still unclear, and confirming whether their university was founded between 1800 and 1825 is essential to fully satisfy the question’s constraints. The PDF linked in the search results appears to be the full text of the original study, which is likely to list the authors’ names and their university affiliations in the header or footer, or in a methods section. By examining this document directly, I can extract the precise university information and determine if it matches the founding date requirement. **I need to visit this PDF now to obtain the definitive author and institutional details.**

</think>

<tool_call>

```
{“name”: “visit”, “arguments”: {“url”: [“https://2024.sci-hub.box/459/fe3e50670a8b274447ead89dd3036a97/10.1016@S0304-40170000274-0.pdf”], “goal”: ”Extract author names and their university affiliations from the PDF”}}
```

</tool_call>

<tool_response>

The useful information in https://2024.sci-hub.box/459/fe3e50670a8b274447ead89dd3036a97/10.1016@S0304-40170000274-0.pdf for user goal Extract author names and their university affiliations from the PDF as follows: ...

</tool_response>

Round 8

<think>

To arrive at the answer, I first **identified** that the study in question was about the jump performance of dog fleas (*Ctenocephalides canis*) and cat fleas (*Ctenocephalides felis felis*), based on search results and abstracts that matched the details in the question—specifically, the use of 450 specimens per species and the linearisation of jump-height curves. I then **confirmed from multiple authoritative sources**, including PubMed and the full-text PDF, that after linearisation, the mean height jump carried out by 50% of fleas was 15.5 cm for *C. canis* and 13.2 cm for *C. felis felis*. Although the university affiliation (Ecole Nationale Vétérinaire de Toulouse) and the recognition period were not explicitly confirmed in the abstracts, the full-text PDF provided the necessary author and institutional details, and the

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

context of the study **matched the timeframes and criteria given in the question**. By cross-referencing these results and verifying the sample sizes and measurement methods, I am confident that the recorded measurements after linearisation are indeed 15.5 cm for *C. canis* and 13.2 cm for *C. felis felis*.
</think>
<answer>
The mean height jump carried out by 50% of the fleas was recorded as 15.5 cm for *Ctenocephalides canis* and 13.2 cm for *Ctenocephalides felis felis*.
</answer> ✓

H USAGE OF LLMs

We employ Qwen2.5-72B-Instruct as the LLM-as-Judge for evaluation, following standard practices in web agent research. This choice is validated by minimal performance variance (< 0.3%) across multiple judging models, including GPT-4o-Mini, Qwen3-235B, Gemini-2.5-Flash. Additionally, we strictly adhere to the official evaluation prompts for each benchmark to ensure comparable assessment.

Apart from the usage of LLMs mentioned in the main text, we additionally use Claude-4 (Anthropic, 2025) to polish paper writing.