

Scaling LLM Agents with Self-Evolving Structured Memory

Anonymous ACL submission

Abstract

With the growing adoption of large language model (LLM) agents in persistent real-world roles, they increasingly face continuous streams of complex tasks requiring iterative reasoning and evidence integration. A common strategy for improving agent performance is inference-time scaling, which allocates additional computation to exploration and reasoning. However, existing self-evolving agents typically scale by expanding interaction trajectories, leaving the knowledge accumulated largely unstructured and making it difficult to consolidate discoveries or guide further reasoning. We address this limitation by reframing inference-time scaling as structured knowledge accumulation rather than trajectory expansion. We propose STRUCTMEM, a framework that represents the agent’s evolving knowledge state as a dynamically constructed knowledge graph. During exploration, newly discovered facts are incrementally integrated into the graph, where structural constraints enable conflict detection, knowledge gap identification, and consistency-aware reasoning. Through this iterative build–verify–expand process, STRUCTMEM progressively constructs a coherent knowledge structure that supports grounded and verifiable answer generation. Experiments on challenging knowledge-intensive benchmarks demonstrate substantial improvements over existing agents in reasoning accuracy and robustness. The code is available at <https://anonymous.4open.science/r/StructMem-code-0104>.

1 Introduction

Large language models (LLMs) are increasingly evolving from standalone question-answering systems into autonomous agents capable of performing complex, long-horizon tasks (Liu et al., 2025; Wang et al., 2024). When deployed in open-world environments, these agents must not only gather information across multiple steps but also continually consolidate newly discovered knowledge and

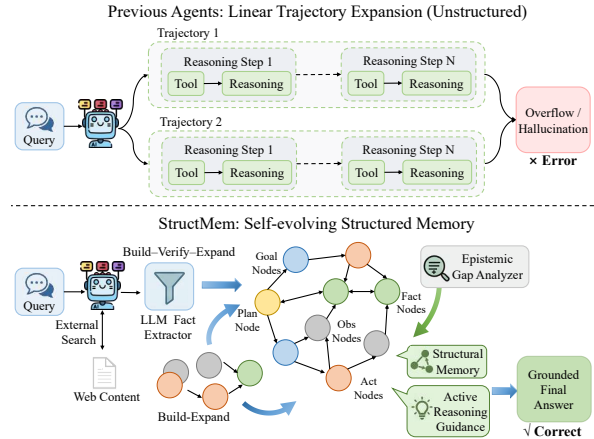


Figure 1: Comparison of inference-time scaling strategies for LLM agents: linear trajectory expansion versus the self-evolving structured memory in STRUCTMEM.

guide subsequent reasoning. Without an effective mechanism to organize intermediate discoveries, the growing context from long interaction histories can easily become noisy and difficult to utilize (Gur et al., 2023; Shi et al., 2025a). To enhance agent capabilities in long-horizon tasks, recent work has explored inference-time scaling, which improves performance by allocating additional computation to reasoning and exploration (Snell et al., 2025). Many self-evolving agents implement this strategy by expanding interaction trajectories or exploring multiple reasoning paths (Zheng et al., 2023b; Kagaya et al., 2024).

Limitations of trajectory-based memory. In these systems, knowledge discovered during exploration is typically accumulated as linear text-based interaction histories or simple memory buffers (Fang et al., 2025). As a result, newly acquired information remains embedded within lengthy reasoning traces rather than being explicitly consolidated into structured knowledge.

This design introduces two fundamental limitations in terms of effectiveness and efficiency. (i) *Effectiveness limitation: lack of explicit reasoning guidance.* Trajectory-based memory buffers primarily serve as passive records rather than

070	structured representations of knowledge. Conse-	generating targeted exploration queries to guide	121
071	quently, they cannot explicitly represent relation-	subsequent reasoning steps. Through this iterative	122
072	ships among retrieved facts, detect logical inconsis-	process, the agent progressively constructs a coherent	123
073	tencies, or reveal missing information required to	knowledge structure, enabling more informed	124
074	solve the task. Without such structural signals, the	exploration and grounded answer generation.	125
075	agent lacks a clear understanding of what knowl-	Experimental results. We evaluate STRUCT-	126
076	edge has been established and what evidence is still	MEM on four challenging deep-search benchmarks:	127
077	required. (ii) <i>Efficiency limitation: lack of effective</i>	BrowseComp (Wei et al., 2025a), BrowseComp-	128
078	<i>information refinement.</i> As agents continuously ap-	ZH (Zhou et al., 2025), GAIA (Mialon et al., 2023),	129
079	pend raw observations to their interaction histories,	and xbench-DeepSearch (Chen et al., 2025). Ex-	130
080	the accumulated context grows rapidly and contains	perimental results show that STRUCTMEM consis-	131
081	substantial redundancy. Without distilling observa-	tently outperforms advanced open-source agents	132
082	tions into compact structured facts, the model must	and proprietary deep-search systems, achieving	133
083	repeatedly process large volumes of unstructured	a significant 19.5% improvement over OpenAI	134
084	text at every reasoning step. This leads to signifi-	DeepResearch (OpenAI, 2025) on the Browse-	135
085	cant computational overhead and poor scalability	Comp benchmark. In addition, experiments on the	136
086	as the interaction trajectory becomes longer.	MedXpertQA benchmark (Zuo et al., 2025) demon-	137
087	Structured memory instead of passive trajecto-	strate strong domain generalization, highlighting	138
088	ries. Insights from cognitive science suggest that	the framework’s ability to support reliable reason-	139
089	human reasoning rarely relies on passive sequential	ing even in highly specialized medical domains	140
090	logs of observations (Bieth et al., 2024). Instead,	without task-specific tuning.	141
091	humans construct structured semantic representa-		
092	tions that organize knowledge into interconnected		
093	relations, enabling them to detect inconsistencies,	2 Related Work	142
094	infer missing information, and guide subsequent		
095	reasoning. This observation suggests that effec-	Deep search agents. Driven by the advanced	143
096	tive long-horizon agent reasoning requires mem-	capabilities of LLMs, recent research has devel-	144
097	ory mechanisms that go beyond flat textual histo-	oped deep search agents to tackle complex, multi-	145
098	ries. As illustrated in Figure 1, rather than pas-	step information-seeking tasks (Nakano et al.,	146
099	sively recording interaction trajectories, an agent’s	2021; Trivedi et al., 2023; Jin et al., 2025; Wei	147
100	working memory should actively organize discov-	et al., 2025b). Unlike basic, single-turn question-	148
101	ered knowledge into a structured representation	answering models, these agents operate in open	149
102	that evolves during exploration.	environments. They use iterative loops of reason-	150
103	Self-evolving structured memory for LLM	ing and tool execution to continuously search the	151
104	agents. Motivated by this perspective, we pro-	web, read pages, and gather evidence (Yao et al.,	152
105	pose STRUCTMEM, a framework that scales LLM	2022; Shinn et al., 2023; Qin et al., 2025; Shi et al.,	153
106	agents through self-evolving structured memory.	2025b). The progress in this field is strongly sup-	154
107	Instead of expanding unstructured interaction tra-	ported by challenging benchmarks like BrowseC-	155
108	jectories, STRUCTMEM represents the agent’s	omp (Wei et al., 2025a) and GAIA (Mialon et al.,	156
109	evolving epistemic state as a dynamically con-	2023), which require agents to persistently navi-	157
110	structed knowledge graph.	gate the internet and synthesize information. To	158
111	During reasoning, the framework performs a con-	succeed on these tasks, modern agents typically	159
112	tinuous build–verify–expand cycle. As the agent	maintain a text-based history of past actions and	160
113	interacts with external environments, raw observa-	observations to plan their future steps (Wang et al.,	161
114	tions are distilled into factual triples by an LLM-	2024; Ouyang et al., 2025).	162
115	based extractor. These facts are integrated into the	Trajectory-based memories face inherent limita-	163
116	existing knowledge graph via a semantic relational	tions. As interaction histories grow longer, early	164
117	integrator, which merges redundant information	information may be forgotten, leading to unstable	165
118	and identifies contradictions. An epistemic gap	behavior such as repeated queries or loss of focus	166
119	analyzer evaluates the evolving graph with respect	on the original goal (Liu et al., 2024). Address-	167
120	to the root task, detecting missing knowledge and	ing these limitations motivates methods that can	168
		organize knowledge more explicitly and guide ex-	169

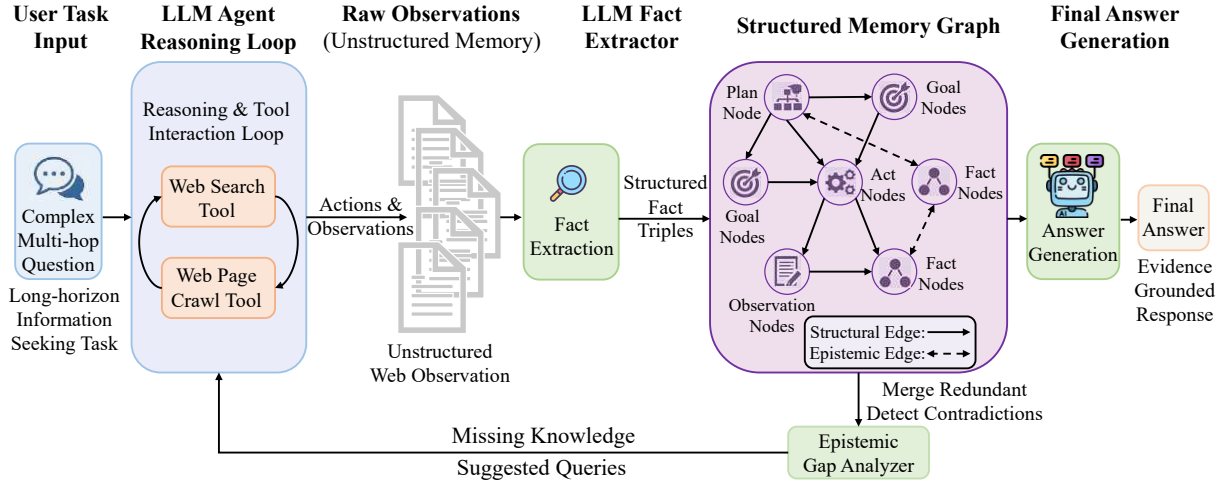


Figure 2: The STRUCTMEM framework for scaling LLM agents. Structured facts are extracted from raw observations, incorporated into a self-evolving knowledge graph, and analyzed for epistemic gaps, enabling guided exploration and graph-grounded answer synthesis.

ploration more effectively.

LLMs for graph reasoning. To enhance reasoning and retrieval, many works integrate LLMs with graph structures. Graph-based retrieval-augmented generation leverages external knowledge graphs to improve multi-hop question answering (Edge et al., 2024; Yasunaga et al., 2022; Zhu et al., 2025). Other approaches, e.g., Tree of Thoughts (Yao et al., 2023) and Graph of Thoughts (Besta et al., 2024), structure the LLM’s internal reasoning process, decomposing complex problems into connected steps and enabling exploration of multiple reasoning paths. Flash-Searcher (Qin et al., 2025) introduces a reasoning architecture based on directed acyclic graphs that enable parallel execution of subtasks and improves accuracy. However, its inference still relies on parallel trajectory expansion, where intermediate knowledge is not explicitly structured and accumulated, leading to redundant exploration and limited reuse of discovered information. This limitation calls for approaches that dynamically construct and evolve structured knowledge during reasoning.

3 Method

To address the efficiency and effectiveness limitations of trajectory-based memory, we propose STRUCTMEM. It integrates a dynamically evolving structured memory directly into the agent’s reasoning loop. STRUCTMEM models the agent’s epistemic state as a knowledge graph that evolves with exploration, via an iterative build–verify–expand cycle (Figure 2) that has three key components: (i) the definition of structural memory elements (§3.1); (ii) dynamic memory evolution via the semantic

relational integrator and epistemic gap analyzer (§3.2); and (iii) graph-grounded answer generation based on the consolidated memory (§3.3).

3.1 Structural memory: definition

We represent the agent’s memory as a directed graph $\mathcal{G}_t = (\mathcal{V}, \mathcal{E})$, dynamically built during task execution. Unlike static knowledge bases, this memory is open-domain and continuously updated by the LLM.

Node types. Nodes represent task-relevant information at varying levels of abstraction: $\mathcal{V} = \{v_{plan}, v_{goal}, v_{action}, v_{obs}, v_{fact}\}$, where

- v_{plan} : high-level strategy parsed from the task.
- v_{goal} : intermediate sub-goals decomposed from the plan.
- v_{action} : executed tool operations (e.g., web search).
- v_{obs} : raw textual feedback retrieved from the environment.
- v_{fact} : structured facts extracted from observations.

To ensure expressiveness, each v_{fact} node is represented as an augmented RDF triples (Klyne, 2004): $\mathcal{F} = \{(e_s, r, e_o, c)\}$. Here, e_s is the subject, r is the relation, e_o is the object, and c is an importance score (e.g., critical or supplementary) assigned by the LLM. The score c helps the system prioritize information when generating the answer.

Edge types. The edge set \mathcal{E} defines the structural and semantic connectivity between nodes. It consists of two main categories: structural dependencies and epistemic relationships.

Structural dependencies track the agent’s execution flow, where (i) the CONTAINS edge repre-

sents general inclusion, such as linking a high-level plan to its specific sub-goals; and (ii) the PRODUCES edge connects a tool action to its resulting raw observation. *Epistemic relationships* map the agreement between different extracted facts, where (i) the SUPPORTS edge indicates corroborating evidence; (ii) the CONTRADICTS edge highlights conflicting information; and (iii) the RELATED edge denotes general semantic relevance.

3.2 Structural memory: evolving

The agent’s memory graph \mathcal{G}_t evolves iteratively from an initial task decomposition, integrating new observations while maintaining global consistency and highlighting knowledge gaps to guide exploration.

Initialization. Given a user task T , an LLM parser decomposes it into a high-level plan and sub-goals, forming the base graph \mathcal{G}_0 . The central v_{plan} node connects to v_{goal} nodes (CONTAINS edges). Task-provided facts are linked to v_{plan} , grounding initial knowledge.

Continuous expansion. At each reasoning step t , a tool action a_t (e.g., a web search or page crawl) produces unstructured observation o_t , prompting a graph update from \mathcal{G}_t to \mathcal{G}_{t+1} . A new v_{action} node represents a_t and links to the governing v_{plan} . The resulting observation o_t is stored in a v_{obs} node connected via a PRODUCES edge. An LLM-based fact extractor $f_{ext}(o_t, T)$ converts o_t into candidate factual triples \mathcal{F}_t , which are instantiated as v_{fact} nodes linked to v_{obs} through CONTAINS edges. This process converts unstructured raw data into structured knowledge while preserving causal relationships.

Fact integration via semantic relational integrator. Before appending new facts to the global graph, a semantic relational integrator evaluates the extracted fact \mathcal{F}_{new} against each existing fact \mathcal{F}_{exist} . Specifically, the LLM computes a pairwise relation $R(\mathcal{F}_{new}, \mathcal{F}_{exist}) \in \{\text{SUPPORTS, CONTRADICTS, RELATED, NONE}\}$. This mechanism actively maintains the memory’s purity: if \mathcal{F}_{new} is semantically identical to an existing fact \mathcal{F}_{exist} , the nodes are merged to prevent redundant context bloat. Conversely, if they conflict, a CONTRADICTS edge is formed. This explicitly highlights an epistemic discrepancy within the memory, signaling the agent to resolve the contradiction in subsequent exploration steps.

Reasoning guidance via epistemic gap analyzer. To prevent the agent from executing redundant queries or failing to recognize task completion, we

introduce an epistemic gap analyzer. The analyzer evaluates the current knowledge against the root plan v_{plan} . Formally, the gap analyzer acts as a deductive function $H(\mathcal{G}_t, T) \rightarrow (\mathcal{K}_{cov}, \mathcal{K}_{mis}, \mathcal{Q}_{sug})$:

- \mathcal{K}_{cov} : The subspace of the task covered by high-confidence factual sub-graphs.
- \mathcal{K}_{mis} : The identified structural information gaps (i.e., missing relationships needed to connect a v_{goal} to v_{fact} nodes).
- \mathcal{Q}_{sug} : A set of precise, dynamically synthesized queries targeting \mathcal{K}_{mis} .

The analysis is prepended to the LLM context, steering subsequent reasoning.

The agent iteratively performs actions, extracts observations, integrates facts, and updates epistemic gaps until either all structural gaps are resolved or a maximum step limit is reached, producing a complete, consistent knowledge graph.

3.3 STRUCTMEM-based answer generation

The agent generates its final response from the consolidated memory graph. Facts are organized by LLM-assigned importance scores, forming a hierarchical context. This structured knowledge, together with the graph’s epistemic confidence, is fed into the final synthesis prompt. Consequently, the LLM acts as an answer generator that explicitly synthesizes its response based on the evidence. This mechanism ensures that all claims in the final answer are grounded in retrieved, structurally verified facts.

4 Experimental Settings

Benchmark datasets. To rigorously validate the information-seeking and reasoning capabilities of STRUCTMEM in complex, human-challenging scenarios, we evaluate our approach on a suite of demanding benchmarks rather than relatively simple question-answering datasets. (i) **BrowseComp** (Wei et al., 2025a): A web browsing benchmark with 1,266 complex tasks designed to evaluate an agent’s ability to perform persistent web browsing and multi-step search to discover hard-to-find information. (ii) **BrowseComp-ZH** (Zhou et al., 2025): A Chinese web browsing benchmark with 289 multi-hop questions that evaluates information-seeking and multihop reasoning in the Chinese web environment. (iii) **GAIA** (Mialon et al., 2023): A benchmark for general AI assistants that tests real-world problem solving requiring reasoning, tool use, and web browsing. We use its text-only validation set. (iv) **xbench-DeepSearch** (Chen et al., 2025): A

subset of the xbench evaluation suite that measures deep-search capability by evaluating search breadth and reasoning depth in complex tasks.

For BrowseComp, we follow (Zeng et al., 2026) and conduct our experiments on the same sampled subset of 100 examples, since the computational costs of running experiments on the full dataset would be too high. The selected 100 examples in BrowseComp were previously subjected to robustness analyses in (Zeng et al., 2026) to ensure their fairness and representativeness. For all other datasets, we conduct evaluations on the full dataset.

Evaluation metrics. We adopt the LLM-as-Judge paradigm (Zheng et al., 2023a; Qin et al., 2025) for automated evaluation, using GPT-4.1-mini as the judge model. For each benchmark, the judge model provides a binary correctness decision for the output generated by the agent. In reporting overall performance, we present accuracy as the default metric.

Baseline models. In our experiments, we compare StructMem with existing methods from three categories. (i) **Proprietary deep research agents:** Grok3 DeepResearch (xAI, 2025a,b), Doubao DeepResearch (ByteDance Doubao, 2025), Tongyi DeepResearch (Tongyi DeepResearch Team, 2025), OpenAI DeepResearch (OpenAI, 2025) and OpenAI ChatGPT agent. (ii) **Advanced LLMs:** GPT-4o (OpenAI, 2024), DeepSeek-R1 (Guo et al., 2025), Gemini-2.5-Pro (Comanici et al., 2025), and OpenAI GPT-5 (Singh et al., 2025). (iii) **Open-source browse agents:** BrowseMaster (Pang et al., 2025), MiroFlow (MiroMind AI Team, 2025) and Flash-Searcher (Qin et al., 2025). As paper that introduced Flash-Searcher does not report results on BrowseComp-ZH, we reproduce its performance using the official implementation under the same evaluation setup. Other model results are taken from prior work (Pang et al., 2025; Zeng et al., 2026; Qin et al., 2025).

Implementation details. We use GPT-5 (Singh et al., 2025) as the base model. It acts as both the main reasoning engine and the manager for all knowledge graph updates. To save computing time and prevent the agent from getting stuck in endless loops, we limit the agent to a maximum of 40 steps per task. The gap analyzer is set to trigger periodically. To make the agent more flexible and avoid fixed patterns, the gap analysis step happens roughly every 8 steps, with slight random changes during the run.

To ensure efficient agent workflows, STRUCTMEM employs a two-tool configuration: (i) A search tool (Serper API; Serper, Inc., 2025): Returns the top 5 relevance-ranked results (title, snippet, URL) per query. This balances comprehensive information access with computational efficiency, preventing information overload. (ii) A crawl tool (Jina Reader; Inc. Jina, 2025): Extract and summarize the online content from a web page. To manage context limits and API costs, pages are truncated to the first 60,000 characters before summarization. This streamlined setup provides robust retrieval capabilities, making it highly effective for the framework’s parallel execution architecture.

5 Experimental Results

We evaluate STRUCTMEM to answer five key research questions: (i) Does structured memory improve reasoning over standalone LLMs and conventional agents? (ii) How does the memory graph adapt to task complexity and compress information? (iii) How sensitive is performance to the choice of backbone LLM? (iv) Can the framework generalize to specialized, knowledge-intensive domains? (v) What do multi-hop reasoning cases reveal about its strengths and weaknesses?

5.1 Main results

Table 1 compares STRUCTMEM with proprietary deep-research agents, advanced reasoning models, and open-source agents across four challenging deep-search benchmarks.

Static models vs. search agents. We first observe a substantial performance gap between standalone reasoning models and agent-based systems. For example, GPT-4o and DeepSeek-R1 achieve only 0.6% and 2.0% accuracy on BrowseComp, respectively. Even the strongest standalone model (OpenAI GPT-5) achieves just 19.8%, which remains noticeably below most agent-based approaches. This trend highlights a key limitation of purely parametric reasoning: complex deep-search tasks require iterative information acquisition and evidence integration, which static models cannot perform.

Performance of existing agents. Agent frameworks significantly improve performance by enabling external search and iterative reasoning. Among the open-source systems, Flash-Searcher achieves the strongest baseline results: 67.7% on BrowseComp and 83.0% on xbench-DeepSearch. Similarly, proprietary systems such as OpenAI DeepResearch demonstrate strong performance on

Model	Benchmark Accuracy (%)			
	BrowseComp	BrowseComp-ZH	GAIA	xbench-DeepSearch
<i>Proprietary Agents</i>				
Grok3 DeepResearch	-	12.9	-	50+
Doubao DeepResearch	-	26.0	-	50+
Tongyi DeepResearch	43.4	46.7	70.9	75.0
OpenAI DeepResearch	51.5	42.9	70.5	66.7
OpenAI ChatGPT agent	68.9	-	-	-
<i>Models</i>				
GPT-4o	0.6	6.2	17.5	18.0
DeepSeek-R1	2.0	23.2	16.5	32.7
Gemini-2.5-Pro	7.6	27.3	-	-
OpenAI GPT-5	19.8	34.3	-	30.0
<i>Open-Source Agents</i>				
BrowseMaster	30.0	46.5	68.0	66.0
MiroFlow	33.2	-	82.4	72.0
Flash-Searcher	67.7	60.2	82.5	83.0
STRUCTMEM	71.0	66.8	83.5	87.0

Table 1: Performance comparison with proprietary deep-research agents, advanced reasoning models, and open-source agents on four challenging deep-search benchmarks.

several benchmarks. These results confirm that search-based reasoning pipelines are crucial for long-horizon information-seeking tasks.

Effectiveness of structured memory. STRUCTMEM consistently achieves the best results across all four benchmarks. In particular, it reaches 71.0% on BrowseComp and 66.8% on BrowseComp-ZH, outperforming the strongest open-source agent Flash-Searcher by +3.3 and +6.6 points respectively. On GAIA and xbench-DeepSearch, STRUCTMEM further improves the best existing results to 83.5% and 87.0%, establishing new state-of-the-art performance.

These improvements highlight the core advantage of STRUCTMEM. While existing agents primarily accumulate reasoning traces as linear trajectories, our framework explicitly organizes discovered information into a dynamically evolving knowledge graph. This structured memory enables more consistent evidence integration and better reasoning guidance, leading to higher accuracy on complex deep-search tasks.

5.2 Structured memory analysis

To examine the internal behavior of STRUCTMEM, we collected statistics from the final knowledge graphs produced by the agent across all tasks in four benchmarks. Figure 3 reports the average number of nodes, edges, v_{obs} nodes, and v_{fact} nodes per task.

Adaptive expansion. As shown in the top panel of Figure 3, on the BrowseComp dataset, the most complex benchmark focusing on deep web explo-

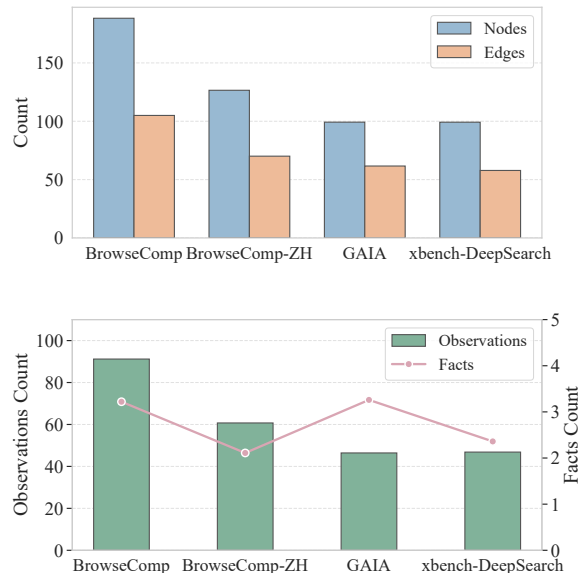


Figure 3: Structured memory statistics. Top: Average nodes vs. edges. Bottom: Raw observations vs. extracted facts.

ration, the graph yields an average of 188.3 nodes and 105.0 edges, nearly twice the scale observed in GAIA (99.2 nodes, 61.7 edges) and xbench-DeepSearch (99.2 nodes, 57.9 edges). This demonstrates that STRUCTMEM adaptively expands in proportion to task complexity: harder queries trigger more tool invocations and denser logical connections, while simpler queries avoid unnecessary computational overhead.

Information compression. The bottom panel of Figure 3 shows that while raw observation nodes are numerous (e.g., 91.2 per task in BrowseComp), the number of integrated fact nodes remains small and stable (averaging only 2–3 per task). This

Model	MedXpertQA Text		
	Reasoning	Understanding	Average
<i>Human Expert (Pre-Licensed)</i>	41.7	45.4	42.6
GPT-4o	30.6	29.5	30.4
GPT-5-nano	36.4	34.0	35.2
GPT-5-mini	45.9	43.8	44.9
GPT-5 (Full Set)	57.0	54.8	55.9
GPT-5 (Sampled Set)	56.0	55.0	55.5
STRUCTMEM	63.0	68.0	65.5

Table 2: Generalization performance on the **MedXpertQA** medical reasoning benchmark.

highlights the effect of the semantic relational integrator: it consolidates observations, removes redundancy, and extracts a compact set of structured facts, yielding a concentrated knowledge representation that filters irrelevant information effectively.

5.3 Backbone model ablation

To disentangle the contribution of STRUCTMEM from the underlying LLM capability, we conduct a backbone ablation by instantiating the framework with models of varying strengths, while keeping all other components unchanged. Experiments are performed on the *xbench-DeepSearch* benchmark.

STRUCTMEM Base Model	<i>xbench-DeepSearch</i>
Gemini-2.5-flash	35.0
GPT-4.1-mini	60.0
GPT-5-mini	78.0
GPT-5 (Default)	87.0

Table 3: Performance of STRUCTMEM with different base models on the *xbench-DeepSearch* benchmark.

As shown in Table 3, we find that: (i) STRUCTMEM performs strongly across all backbone models, including lightweight LLMs, indicating that its effectiveness does not rely on large-scale parametric knowledge; (ii) performance improves steadily with model capability (from 35.0 to 87.0), showing that STRUCTMEM scales synergistically with stronger backbones; and (iii) the strong performance of smaller models (e.g., GPT-5-mini reaching 78.0) highlights that STRUCTMEM provides a robust external reasoning substrate that compensates for limited model capacity. This confirms that the gains stem from structured reasoning and memory organization, rather than solely from the backbone LLM. Overall, these results demonstrate that STRUCTMEM is a backbone-agnostic framework that consistently enhances reasoning performance, offering both strong effectiveness with small models and scalable gains with larger ones.

5.4 Model generalization analysis

To assess the cross-domain generalization ability of STRUCTMEM, we evaluate the framework on the **MedXpertQA** benchmark (Zuo et al., 2025), which requires expert-level medical reasoning and deep understanding (Wang et al., 2025). We focus on the text-only version, consisting of two question types: Reasoning and Understanding. Due to computational constraints, we randomly sampled 100 instances from each subset. To confirm representativeness, we benchmarked the standard GPT-5 model on this sample, obtaining 56.0% on Reasoning and 55.0% on Understanding, closely matching full-set results (57.0% and 54.8%). This validates that the subset provides a reliable basis for comparison.

Table 2 reports the results. STRUCTMEM achieves the highest overall average score (65.5%), surpassing both GPT-5 (55.9%) and the pre-licensed human expert baseline (42.6%). This demonstrates that the framework generalizes effectively beyond its original web-search domain, supporting reasoning over specialized, knowledge-intensive content. The largest improvement is observed on the Understanding subset (+13.2% over GPT-5), suggesting that tasks requiring synthesis of complex concepts benefit disproportionately from STRUCTMEM’s dynamic structured memory, which organizes scattered evidence into coherent, relational representations. Overall, these results confirm that dynamically structuring knowledge not only enhances performance in open-world environments but also translates into robust domain generalization for specialized reasoning tasks.

5.5 Case study

We present a representative multi-hop reasoning trajectory from BrowseComp (Task ID: 28) in Table 4 to provide a qualitative analysis of STRUCTMEM. **Strengths.** This case illustrates how STRUCTMEM transforms large-scale, noisy web interactions into

Component	Content detail
Complex task (Multi-hop query)	A TV show in the 1960s was produced by a TV station, channel 8, which, as of 2023, ran programming from the PBS network, in a city in the US that shares its name with a city in Southern Europe. The show’s presenter interacted with cartoon characters, and the name of the show includes a food item. In the series, the presenter’s first name was only two letters. The show had an episode whose title mentioned a trip to a specific destination. What was this destination?
Initial plan & goal (Nodes: 3)	Initial plan: Targeted PBS station verification. Use web search with queries like “WGTV channel 8 Athens PBS” to identify WGTV as a PBS member station located in Athens, GA (sharing a name with Athens, Greece)... Goal 1: Identify the qualifying PBS “channel 8” station and its city that shares a name with a city in Southern Europe... Goal 2: ...
Raw observations (Nodes: 124) [Highly noisy]	<i>The agent executed 124 search/crawl actions and accumulated massive unstructured text buffers:</i> <ul style="list-style-type: none"> • Obs_2: [WGTV CH 8 Georgia Public Broadcasting] (https://www.gpb.org/...) • Obs_4: [Georgia Greats Dean Rusk: At the Heartbeat of History - PBS] ... • Obs_24: [PDF] ED119715 - Produced by WGTV-TV, Athens GA; episode list; A Trip to the University; series title; host... • Obs_37: TV Guide Archive - PBS Schedule 1974; WGTV programming lineup; evening news; documentary slot; local broadcast schedule... <i>... (121 other lengthy webpage records, historical archives, and irrelevant TV guides omitted)</i>
Structured facts (Nodes: 5) [Significant compression]	STRUCTMEM’s <i>Semantic Relational Integrator</i> distilled the 124 noisy web pages into exactly 5 high-value factual triples: <ol style="list-style-type: none"> 1. ⟨Show episode, title, “A Trip to the University”⟩ 2. ⟨“A Trip to the University”, destination, “the University”⟩ 3. ⟨Show episode 21, title, “Visit to the World of the Muffins, Part I”⟩ 4. ⟨Miss Jo and her cartoon drawings, <i>take children to</i>, the curious world of Muffin land⟩ 5. ⟨The show, <i>encourages</i>, 3 to 7-year-olds to use imagination⟩
Agent’s prediction	the University
Standard answer	University (Correct match)

Table 4: Case study from the BrowseComp benchmark. STRUCTMEM successfully prevents context window overflow by distilling 124 unstructured and noisy webpage observations into just 5 explicit factual triples. These structured facts enable the agent to accurately deduce the standard answer to a highly complex multi-hop query. Due to space limitations, we only show several key sentences from the reasoning process.

557 concise and actionable structured knowledge. Start-
558 ing from 124 raw observations collected via search
559 and browsing, the framework distills only a small
560 set of key fact triples through fact extraction and
561 relational integration. This structured memory ex-
562 plicitly captures core entities and their relations,
563 enabling the agent to perform stable multi-hop rea-
564 soning without being distracted by irrelevant con-
565 text. Compared to trajectory-based methods over
566 unstructured histories, this yields better informa-
567 tion compression and reasoning focus.

568 **Weaknesses.** Noisy observations may still be ab-
569 stracted into spurious fact triples, which can propa-
570 gate and affect downstream reasoning. Thus, while
571 structured memory improves robustness, it does not
572 fully eliminate error accumulation, especially when
573 early extraction steps are inaccurate. This suggests
574 the need for more reliable fact validation and noise
575 filtering mechanisms (see Section Limitations).

576 Overall, STRUCTMEM effectively condenses
577 large-scale, noisy web observations into concise
578 structured knowledge for accurate reasoning, while

leaving robustness as a key area for improvement. 579

6 Conclusion 580

581 We have introduced STRUCTMEM, a novel agent
582 framework for structured reasoning on complex,
583 long-horizon tasks. Unlike conventional agents
584 that accumulate unstructured interaction traces,
585 STRUCTMEM maintains a compact, dynamically
586 evolving memory by constructing knowledge
587 graphs on-the-fly. The semantic relational integra-
588 tor distills and organizes noisy observations into a
589 concise set of structured facts, while the epistemic
590 gap analyzer identifies unresolved knowledge gaps,
591 generates targeted exploratory queries, and guides
592 the agent’s reasoning in real time.

593 Extensive experiments show that STRUCTMEM
594 achieves state-of-the-art performance across deep-
595 search benchmarks and generalizes to specialized
596 domains without task-specific fine-tuning. Overall,
597 STRUCTMEM demonstrates that structured mem-
598 ory can serve as a scalable and effective abstraction
599 layer for inference-time reasoning.

600 Limitations

601 Despite its strong performance, STRUCTMEM in-
602 troduces several limitations. (i) *Error propagation*
603 *in structured abstraction*. The framework relies
604 on the underlying LLM to extract and structure
605 facts from noisy observations. When the backbone
606 model produces inaccurate or incomplete abstrac-
607 tions, these errors can be explicitly encoded into the
608 structured memory and propagate through subse-
609 quent reasoning steps. While structuring improves
610 interpretability and efficiency, it also makes errors
611 more persistent once incorporated into the memory.
612 (ii) *Efficiency–effectiveness trade-off*. The itera-
613 tive build–verify–expand process requires multiple
614 LLM calls for fact extraction, relational integra-
615 tion, and gap analysis, leading to increased com-
616 putational overhead and inference latency. This
617 highlights an inherent trade-off between structured
618 reasoning quality and efficiency.

619 Future work will explore more robust and ef-
620 ficient structured reasoning mechanisms. On the
621 robustness side, we plan to incorporate uncertainty-
622 aware fact representations (e.g., confidence-
623 weighted or probabilistic edges) and develop verifi-
624 cation strategies based on cross-source agreement
625 or multi-path reasoning consistency, enabling the
626 agent to detect and revise unreliable facts dynam-
627 ically. On the efficiency side, we aim to distill
628 the structured memory construction process into
629 lightweight, specialized modules, and investigate
630 adaptive computation strategies (e.g., selectively
631 triggering memory updates or reasoning steps) to
632 reduce unnecessary overhead. More broadly, an
633 interesting direction is to enable self-refining mem-
634 ory, where the agent continuously revises and reor-
635 ganizes its structured knowledge over time, poten-
636 tially improving both accuracy and long-horizon
637 reasoning efficiency.

638 Ethical Considerations

639 In this paper, all the foundation models and tools
640 used in our experiments are publicly accessible.
641 All experimental data and evaluations are based
642 on publicly available benchmarks, and no new
643 personal information is exposed or collected in
644 this process, ensuring that user privacy is strictly
645 protected. Additionally, the STRUCTMEM frame-
646 work we propose aims to enhance the effectiveness,
647 transparency, and factual grounding of autonomous
648 agents during complex reasoning tasks. By ex-
649 plicitly structurally resolving factual conflicts, our

approach actively mitigates LLM hallucinations
and does not encourage or induce the model to
produce any harmful, biased, or misleading infor-
mation. Therefore, We believe our work complies
with ACL ethical guidelines.

References

- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gersten-
berger, Michal Podstawski, Lukas Gianinazzi, Joanna
Gajda, Tomasz Lehmann, Hubert Niewiadomski, Pi-
otr Nyczyk, and 1 others. 2024. Graph of thoughts:
Solving elaborate problems with large language mod-
els. In *Proceedings of the AAAI conference on artifi-
cial intelligence*, volume 38, pages 17682–17690.
- Théophile Bieth, Yoed N Kenett, Marcela Ovando-
Tellez, Alizée Lopez-Persem, Célia Lacaux, Marie
Scuccimarra, Inès Maye, Jade Sénéchal, Delphine
Oudiette, and Emmanuelle Volle. 2024. Changes
in semantic memory structure support successful
problem-solving and analogical transfer. *Communi-
cations Psychology*, 2(1):54.
- ByteDance Doubao. 2025. Doubao. <http://www.doubao.com/>.
- Kaiyuan Chen, Yixin Ren, Yang Liu, Xiaobo Hu, Hao-
tong Tian, Tianbao Xie, Fangfu Liu, Haoye Zhang,
Hongzhang Liu, Yuan Gong, and 1 others. 2025.
xbench: Tracking agents productivity scaling with
profession-aligned real-world evaluations. *arXiv
preprint arXiv:2506.13651*.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann,
Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-
cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and
1 others. 2025. Gemini 2.5: Pushing the frontier with
advanced reasoning, multimodality, long context, and
next generation agentic capabilities. *arXiv preprint
arXiv:2507.06261*.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua
Bradley, Alex Chao, Apurva Mody, Steven Truitt,
Dasha Metropolitanaky, Robert Osazuwa Ness, and
Jonathan Larson. 2024. From local to global: A
graph rag approach to query-focused summarization.
arXiv preprint arXiv:2404.16130.
- Runnan Fang, Yuan Liang, Xiaobin Wang, Jialong
Wu, Shuofei Qiao, Pengjun Xie, Fei Huang, Hua-
jun Chen, and Ningyu Zhang. 2025. Memp: Ex-
ploring agent procedural memory. *arXiv preprint
arXiv:2508.06433*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao
Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu
Zhang, Shirong Ma, Xiao Bi, and 1 others. 2025.
Deepseek-r1: Incentivizing reasoning capability in
llms via reinforcement learning. *arXiv preprint
arXiv:2501.12948*.

702	Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. 2023. A real-world webagent with planning, long context understanding, and program synthesis. <i>arXiv preprint arXiv:2307.12856</i> .	755
703		756
704		757
705		758
706		759
707	Inc. Jina. 2025. Jina reader. https://jina.ai/reader/ .	760
708		761
709	Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. <i>arXiv preprint arXiv:2503.09516</i> .	762
710		763
711		764
712		765
713		766
714	Tomoyuki Kagaya, Thong Jing Yuan, Yuxuan Lou, Jayashree Karlekar, Sugiri Pranata, Akira Kinose, Koki Oguri, Felix Wick, and Yang You. 2024. Rap: Retrieval-augmented planning with contextual memory for multimodal llm agents. <i>arXiv preprint arXiv:2402.03610</i> .	767
715		768
716		769
717		770
718		771
719		772
720	Graham Klyne. 2004. Resource description framework (rdf): Concepts and abstract syntax. http://www.w3.org/TR/rdf-concepts/ .	773
721		774
722		775
723	Bang Liu, Xinfeng Li, Jiayi Zhang, Jinlin Wang, Tanjin He, Sirui Hong, Hongzhang Liu, Shaokun Zhang, Kaitao Song, Kunlun Zhu, and 1 others. 2025. Advances and challenges in foundation agents: From brain-inspired intelligence to evolutionary, collaborative, and safe systems. <i>arXiv preprint arXiv:2504.01990</i> .	776
724		777
725		778
726		779
727		780
728		781
729		782
730	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. <i>Transactions of the association for computational linguistics</i> , 12:157–173.	783
731		784
732		785
733		786
734		787
735	Grégoire Mialon, Clémentine Fourier, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023. Gaia: a benchmark for general ai assistants. In <i>The Twelfth International Conference on Learning Representations</i> .	788
736		789
737		790
738		791
739		792
740	MiroMind AI Team. 2025. Miroflow: An open-source agentic framework for deep research. https://github.com/MiroMindAI/MiroFlow . GitHub repository.	793
741		794
742		795
743		796
744	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, and 1 others. 2021. Webgpt: Browser-assisted question-answering with human feedback. <i>arXiv preprint arXiv:2112.09332</i> .	797
745		798
746		799
747		800
748		801
749		802
750	OpenAI. 2024. Hello gpt-4o. https://openai.com/index/hello-gpt-4o/ .	803
751		804
752		805
753	OpenAI. 2025. Deep research system card. https://cdn.openai.com/deep-research-system-card.pdf .	806
754		807
		808
		809
		810
	Siru Ouyang, Jun Yan, I Hsu, Yanfei Chen, Ke Jiang, Zifeng Wang, Rujun Han, Long T Le, Samira Daruki, Xiangru Tang, and 1 others. 2025. Reasoningbank: Scaling agent self-evolving with reasoning memory. <i>arXiv preprint arXiv:2509.25140</i> .	
	Xianghe Pang, Shuo Tang, Rui Ye, Yuwen Du, Yaxin Du, and Siheng Chen. 2025. Browsemaster: Towards scalable web browsing via tool-augmented programmatic agent pair. <i>arXiv preprint arXiv:2508.09129</i> .	
	Tianrui Qin, Qianben Chen, Sinuo Wang, He Xing, King Zhu, He Zhu, Dingfeng Shi, Xinxin Liu, Ge Zhang, Jiaheng Liu, and 1 others. 2025. Flash-searcher: Fast and effective web agents via dag-based parallel execution. <i>arXiv preprint arXiv:2509.25301</i> .	
	Serper, Inc. 2025. Serper api. https://serper.dev/ .	
	Zhengliang Shi, Yiqun Chen, Haitao Li, Weiwei Sun, Shiyu Ni, Yougang Lyu, Run-Ze Fan, Bowen Jin, Yixuan Weng, Minjun Zhu, Qiuji Xie, Xinyu Guo, Qu Yang, Jiayi Wu, Jujia Zhao, Xiaqiang Tang, Xinbei Ma, Cunxiang Wang, Jiaxin Mao, and 7 others. 2025a. Deep research: A systematic survey. Preprint , arXiv:2512.02038.	
	Zhengliang Shi, Lingyong Yan, Dawei Yin, Suzan Verberne, Maarten de Rijke, and Zhaochun Ren. 2025b. Iterative self-incentivization empowers large language models as agentic searchers. <i>arXiv preprint arXiv:2505.20128</i> .	
	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. <i>Advances in neural information processing systems</i> , 36:8634–8652.	
	Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, Adi Ganesh, Ahmed El-Kishky, Aidan McLaughlin, Aiden Low, AJ Ostrow, Akhila Ananthram, and 1 others. 2025. Openai gpt-5 system card. <i>arXiv preprint arXiv:2601.03267</i> .	
	Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2025. Scaling llm test-time compute optimally can be more effective than scaling parameters for reasoning. In <i>The Thirteenth International Conference on Learning Representations</i> .	
	Tongyi DeepResearch Team. 2025. Tongyi-deepresearch . GitHub repository.	
	Harsh Trivedi, Niranjana Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. In <i>Proceedings of the 61st annual meeting of the association for computational linguistics (volume 1: long papers)</i> , pages 10014–10037.	
	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, and 1 others. 2024. A survey on large language model based autonomous agents. <i>Frontiers of Computer Science</i> , 18(6):186345.	

811	Shansong Wang, Mingzhe Hu, Qiang Li, Mojtaba Safari, and Xiaofeng Yang. 2025. Capabilities of gpt-5 on multimodal medical reasoning. <i>arXiv preprint arXiv:2508.08224</i> .	Browsecomp-zh: Benchmarking web browsing ability of large language models in chinese. <i>arXiv preprint arXiv:2504.19314</i> .	865
812			866
813			867
814			
815	Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. 2025a. Browsecomp: A simple yet challenging benchmark for browsing agents. <i>arXiv preprint arXiv:2504.12516</i> .	Yihua Zhu, Qianying Liu, Akiko Aizawa, and Hidetoshi Shimodaira. 2025. Beyond chains: Bridging large language models and knowledge bases in complex question answering. <i>arXiv preprint arXiv:2505.14099</i> .	868
816			869
817			870
818			871
819			872
820			
821	Wenda Wei, Yu-An Liu, Ruqing Zhang, Jiafeng Guo, Lixin Su, Shuaiqiang Wang, Dawei Yin, Maarten de Rijke, and Xueqi Cheng. 2025b. Thinking forward and backward: Multi-objective reinforcement learning for retrieval-augmented reasoning. <i>arXiv preprint arXiv:2511.09109</i> .	Yuxin Zuo, Shang Qu, Yifei Li, Zhangren Chen, Xuekai Zhu, Ermo Hua, Kaiyan Zhang, Ning Ding, and Bowen Zhou. 2025. Medxpertqa: Benchmarking expert-level medical reasoning and understanding. <i>arXiv preprint arXiv:2501.18362</i> .	873
822			874
823			875
824			876
825			877
826			
827	xAI. 2025a. Grok 3 beta — the age of reasoning agents. https://x.ai/news/grok-3 .		
828			
829	xAI. 2025b. Grok 4 — the most intelligent model in the world. https://x.ai/news/grok-4 .		
830			
831	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. <i>Advances in neural information processing systems</i> , 36:11809–11822.		
832			
833			
834			
835			
836	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. In <i>The eleventh international conference on learning representations</i> .		
837			
838			
839			
840			
841	Michihiro Yasunaga, Antoine Bosselut, Hongyu Ren, Xikun Zhang, Christopher D Manning, Percy S Liang, and Jure Leskovec. 2022. Deep bidirectional language-knowledge graph pretraining. <i>Advances in Neural Information Processing Systems</i> , 35:37309–37323.		
842			
843			
844			
845			
846			
847	Weihao Zeng, Keqing He, Chuqiao Kuang, Xiaoguang Li, and Junxian He. 2026. Pushing test-time scaling limits of deep search with asymmetric verification. In <i>The Fourteenth International Conference on Learning Representations</i> .		
848			
849			
850			
851			
852	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023a. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in neural information processing systems</i> , 36:46595–46623.		
853			
854			
855			
856			
857			
858	Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. 2023b. Synapse: Trajectory-as-exemplar prompting with memory for computer control. <i>arXiv preprint arXiv:2306.07863</i> .		
859			
860			
861			
862	Peilin Zhou, Bruce Leon, Xiang Ying, Can Zhang, Yifan Shao, Qichen Ye, Dading Chong, Zhiling Jin, Chenxuan Xie, Meng Cao, and 1 others. 2025.		
863			
864			