

Evo-Memory: Benchmarking LLM Agent Test-time Learning with Self-Evolving Memory

Tianxin Wei^{†,1}, Noveen Sachdeva², Benjamin Coleman², Zhankui He², Yuanchen Bei¹, Xuying Ning¹, Mengting Ai¹, Yunzhe Li^{†,1}, Jingrui He¹, Ed H. Chi², Chi Wang², Shuo Chen², Fernando Pereira², Wang-Cheng Kang² and Derek Zhiyuan Cheng²

[†]Work done while at Google DeepMind, ¹University of Illinois Urbana-Champaign, ²Google DeepMind

Statefulness is essential for large language model (LLM) agents to perform long-term planning and problem-solving. This makes *memory* a critical component, yet its management and evolution remain largely underexplored. Existing evaluations mostly focus on static conversational settings, where memory is passively retrieved from dialogue to answer queries, overlooking the dynamic ability to accumulate and reuse *experience* across evolving task streams. In real-world environments such as interactive problem assistants or embodied agents, LLMs are required to handle continuous task streams, yet often fail to learn from accumulated interactions, losing valuable contextual insights, a limitation that calls for *test-time evolution*, where LLMs retrieve, integrate, and update memory continuously during deployment. To bridge this gap, we introduce Evo-Memory, a comprehensive streaming benchmark and framework for evaluating *self-evolving memory* in LLM agents. Evo-Memory structures datasets into sequential task streams, requiring LLMs to search, adapt, and evolve memory after each interaction. We unify and implement over ten representative memory modules and evaluate them across 10 diverse multi-turn goal-oriented and single-turn reasoning and QA datasets. To better benchmark experience reuse, we provide a baseline method, ExpRAG, for retrieving and utilizing prior experience, and further propose ReMem, an *action-think-memory refine* pipeline that tightly integrates reasoning, task actions, and memory updates to achieve continual improvement.

Keywords: LLMs, Agentic Memory, Test-time Learning, Self-evolving Agents, Lifelong Intelligence

1. Introduction

Large Language Models (LLMs) have rapidly evolved from simple chatbots into capable systems that can write code, control browsers, and perform advanced question answering (Comanici et al., 2025). These advances have been driven by improving inference, planning, and tool use, as shown by benchmarks emphasizing logical reasoning and multi-step actions. Yet a fundamental capability, *memory*, remains largely underexplored. Memory allows LLMs to maintain state across interactions, accumulate experience, and adapt strategies over time. Recent studies have introduced memory modules that track dialogue histories through compression, indexing, or retrieval (Maharana et al., 2024b), improving *conversational recall* and personalization. However, most of these systems only reuse static dialogue context rather than learning from experience to improve future reasoning or decision-making.

Despite these advances, existing LLM memory systems remain largely static, retrieving information passively rather than evolving through use. Current evaluations test whether models can recall past context but rarely assess their ability to *reuse experience*. In essence, agents remember what was said but not what was learned. *Conversational recall* retrieves prior facts, whereas *experience reuse* abstracts reasoning strategies for future tasks. Without such reuse, models repeatedly solve similar problems, as long-term assistants often recall context yet fail to adapt across sessions.

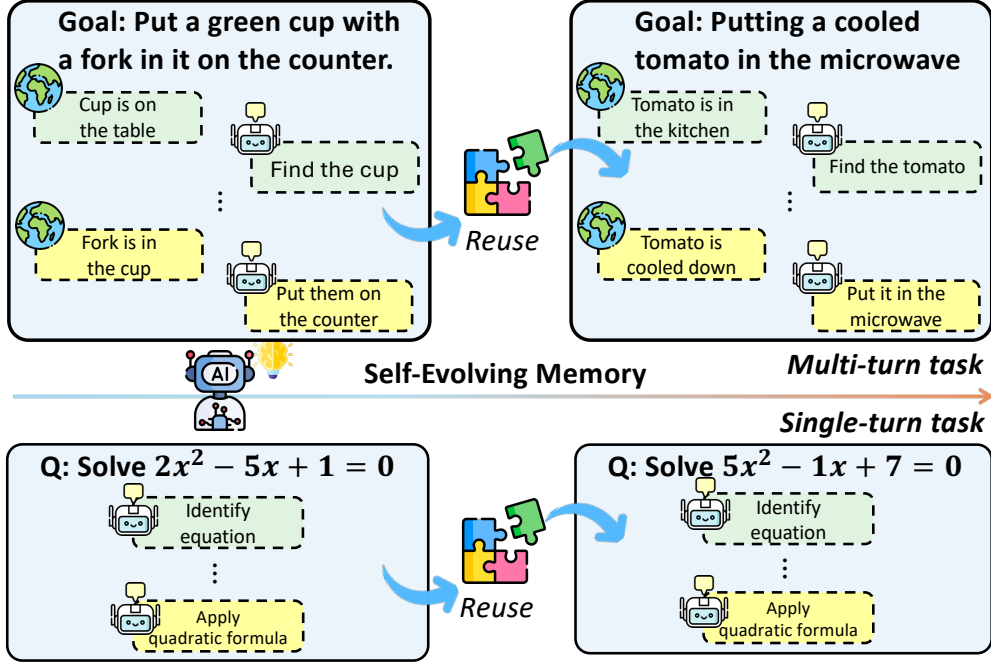


Figure 2 | Illustration of different task types and experience reusing. A stateful agent encounters both multi-turn tasks (e.g., embodied manipulation) and single-turn tasks (e.g., solving equations), and should learn reusable experiences from past experiences.

Several recent benchmarks have begun examining static adaptation but remain limited in scope. StreamBench (Wu et al., 2024a) evaluates sequential learning but mainly measures factual retention without reasoning or trajectory reuse. Lifelong-Bench (Zheng et al., 2025) studies lifelong learning across environments and skills but focuses on retention without modeling memory structure or updates. Other studies (Hu et al., 2025; Maharana et al., 2024b; Wu et al.) assess long-term conversational consistency but do not test how agents evolve their memory during deployment. Together, these efforts highlight a critical gap: while progress has been made on sequential reasoning, there is still no unified framework for evaluating how different memory methods retrieve, integrate, and evolve historical strategies in realistic streaming scenarios. Figure 1 illustrates this contrast between static recall and cumulative improvement through self-evolving memory.

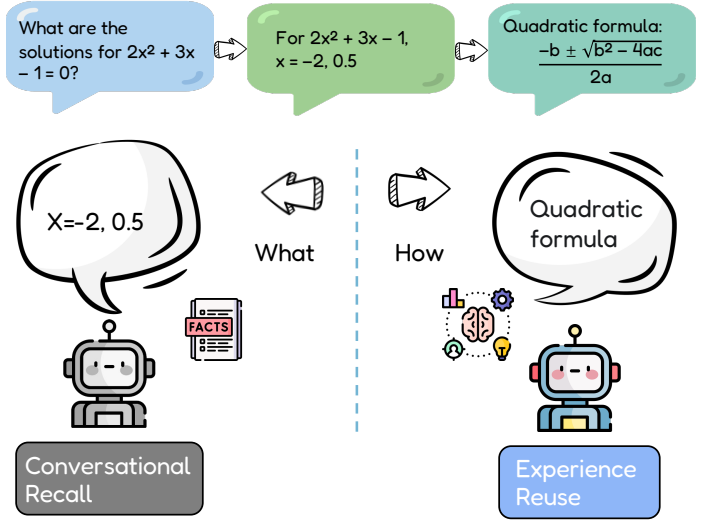


Figure 1 | Conversational recall retrieves past facts (e.g., solutions to $2x^2 + 3x - 1 = 0$). Experience reuse recalls reasoning strategies (e.g., using the formula).

To bridge this gap, we introduce **Evo-Memory**, a comprehensive streaming benchmark and framework for evaluating *self-evolving memory* in LLM agents. Figure 2 illustrates how a self-evolving agent reuses prior experiences across both multi-turn interactive tasks and single-turn reasoning

tasks. Evo-Memory restructures datasets into sequential *task streams*, requiring models to retrieve, adapt, and evolve memory after each interaction. The benchmark covers both *multi-turn goal-oriented* environments and *single-turn reasoning or problem-solving* tasks, explicitly testing whether LLMs can accumulate knowledge and refine strategies during deployment, a process we term *test-time evolution*. We unify and implement over ten representative memory modules, including retrieval-based, workflow, and hierarchical memory systems, to study their adaptation behavior. To further examine experience reuse, we introduce **ExpRAG**, a simple retrieval-based baseline that leverages prior task experiences, and further develop **ReMem**, an advanced *action–think–memory refine* pipeline that tightly integrates reasoning, action, and memory updates for continual improvement.

In summary, our contributions are threefold:

- **Benchmark:** We present Evo-Memory, a streaming benchmark that evaluates LLM agents’ ability to perform *test-time evolution* across diverse multi-turn and single-turn tasks, bridging the gap between conversational recall and experience reuse.
- **Framework:** We provide a unified evaluation framework with memory-centric metrics for analyzing adaptation, efficiency, and stability, and will release all code and configurations for reproducibility.
- **Analysis and Insights:** We introduce **ExpRAG**, a simple retrieval-based baseline for experience reuse, and **ReMem**, an *action–think–memory refine* pipeline that unifies reasoning, action, and memory for continual improvement, informing future designs of memory.

2. Related Work

In this section, we review existing works on test-time learning and self-evolving memory.

2.1. Test-time Learning

Test-time learning (TTL) builds upon early work on test-time adaptation (TTA) (Niu et al., 2022; Wang et al., 2021; Zhang et al., 2023), which enables models to adjust to distribution shifts during deployment. Recent advances extend TTA toward *continuous self-improvement* (Iwasawa and Matsuo, 2021; Liu et al., 2023), allowing models to refine their behavior through online optimization. Recent *agent-based* studies operationalize such continual improvement via reflection, planning, and self-evolution. Works like (Park et al., 2023; Shinn et al., 2023; Wang et al., 2023; Zhao et al., 2025; Zhou et al., 2024) and newer frameworks, including (Chen et al., 2025; Huang et al., 2025) demonstrate how agents autonomously revise plans, synthesize feedback, and co-evolve (Gao et al., 2025). These advances mark a shift from static adaptation toward adaptive, self-improving agents capable of continual learning during deployment. Building on this trend, we propose to benchmark such dynamics from a novel *self-evolving memory* perspective.

2.2. Self-evolving Memory

Early LLM memory systems primarily served as *passive storage*, maintaining recent dialogues or retrieved facts to compensate for limited context windows (Asai et al., 2024a; Lewis et al., 2020; Liu, 2022; Packer et al., 2023; Zhong et al., 2023). Subsequent studies introduced richer management mechanisms, including differentiable read–write controllers (Liang et al., 2023; Modarressi et al., 2023) and evaluations under realistic conversational settings (Maharana et al., 2024a; Wu et al., 2024b). Beyond static buffers, recent work explores *policy-driven control*, where the model is explicitly optimized to decide what to store, retrieve, or overwrite (Li et al., 2025; Xu et al., 2025; Yan et al., 2025; Yu et al., 2025; Zhou et al., 2025). Meanwhile, structured memory representations have emerged to organize experiences into relational or procedural forms, as in RepoGraph (Ouyang et al., 2024),

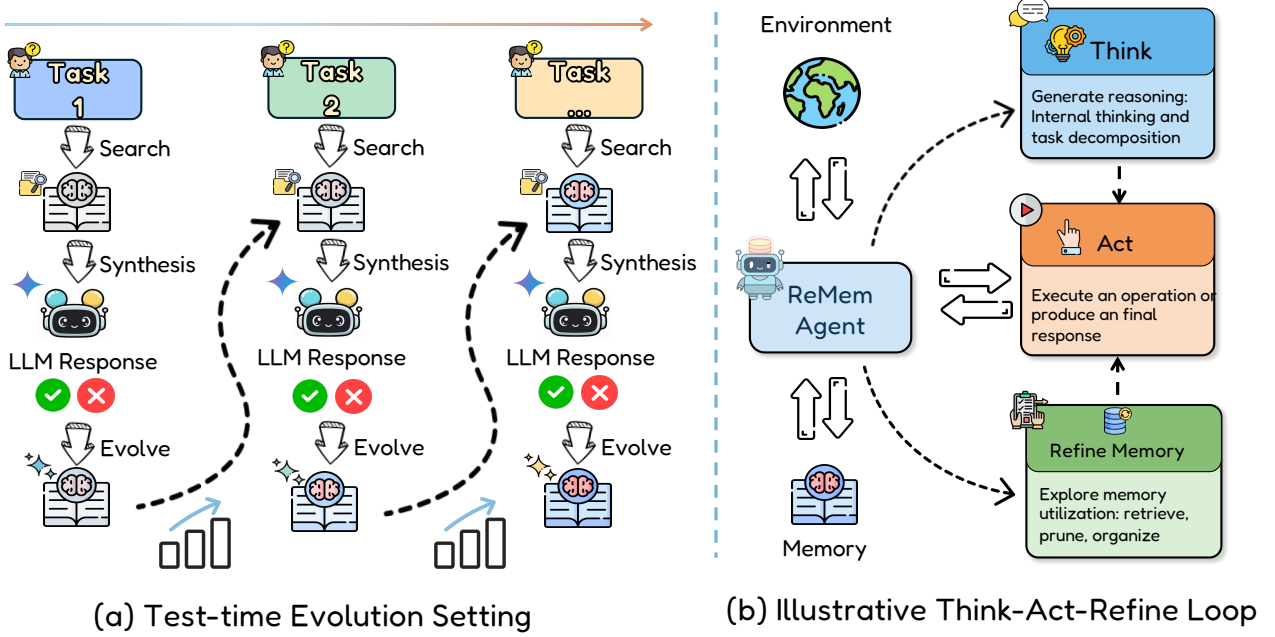


Figure 3 | Overview of the ReMem agent framework. Left: Test-time evolution process where the agent iteratively searches, synthesizes, and evolves its memory across multiple tasks. Right: Agent architecture with three core modules—Think (reasoning and decomposition), Refine Memory (retrieve, prune, organize), and Act (execution)—that interact with the environment and learned memory.

MEM0 (Chhikara et al., 2025), Zep (Rasmussen et al., 2025), and Dynamic Cheatsheets (Suzgun et al., 2025). However, there remains no unified evaluation setting and framework for *self-evolving memory*, the ability to reuse and adapt experiences across tasks. Evo-Memory builds on this trajectory by benchmarking how LLMs not only store and recall but also evolve, reorganize, and reuse memory under streaming task settings.

3. Evo-Memory: Evaluating Self-Evolving Memory in LLM Agents

Existing evaluations of LLMs often treat memory as static recall, overlooking its role in continual adaptation. Evo-Memory provides a unified benchmark to study *self-evolving memory*, where agents retrieve, integrate, and update knowledge over time. As illustrated in Figure 3, the left side shows the test-time evolution process, and the right side outlines the ReMem agent with three modules—Think, Act, and Refine Memory. We first formalize the problem setting, then describe two representative implementations, EXP-RAG and ReMem, used to instantiate the benchmark.

3.1. Problem Formulation

We formalize a general memory-augmented agent as a tuple (F, U, R, C) , where F is the base LLM, U is the memory update pipeline, R is the retrieval module, and C is the contextual construction mechanism that transforms retrieved content into the final working context. In our setting, the agent processes a sequence of inputs $\{x_1, x_2, \dots, x_T\}$, and the memory state M_t evolves with the history. At time t , the agent receives an input x_t , maintains an evolving memory M_t , retrieves relevant elements $R(M_t, x_t)$, constructs a contextualized prompt

$$C_t = C(x_t, R(M_t, x_t)),$$

and produces an output

$$\hat{y}_t = F(C_t).$$

This abstraction unifies a wide spectrum of existing memory mechanisms, from retrieval-augmented generation to dynamic, hierarchical, and workflow-based memories, under a single iterative formulation.

Search. Given the current input x_t , the agent first retrieves relevant memory entries:

$$R_t = R(M_t, x_t),$$

where R can represent similarity search, index-based lookup, or attention over stored embeddings. This step captures memory access policies across different algorithms.

Synthesis. The agent interprets and restructures the retrieved information R_t into a concise working context aligned with the current input x_t . This synthesis yields a coherent text \tilde{C}_t , from which the final output is derived:

$$\hat{y}_t = F(\tilde{C}_t).$$

Synthesis. The agent restructures the retrieved information R_t into a working context tailored to the current input x_t . This step may involve forming a structured prompt (Wang et al., 2024), selecting key memory items (Chhikara et al., 2025; Xu et al., 2025), or merging retrieved content (Suzgun et al., 2025) into a short summary. We denote the resulting context as $\tilde{C}_t = C(x_t, R_t)$, and the final output is

$$\hat{y}_t = F(\tilde{C}_t).$$

Evolve. After obtaining \hat{y}_t , the agent constructs a new memory entry $m_t = h(x_t, \hat{y}_t, f_t)$ that captures the current step’s experience together with the feedback f_t , such as whether the task was completed. The memory is then updated via:

$$M_{t+1} = U(M_t, m_t).$$

Different algorithms instantiate U differently, for example, direct append for retrieval-based memories, summarization or compression for long-term storage, or replacement for bounded-capacity stores. This unified formulation abstracts the essential cycle of *retrieval*, *synthesis*, and *evolution* underlying all memory-based agents.

Dataset Preparation. Evo-Memory restructures conventional static datasets into *streaming task sequences*, enabling evaluation of how LLMs reuse and evolve memory over time. Each dataset can thus be transformed into a sequence $\tau = \{(x_1, y_1), \dots, (x_T, y_T)\}$, forming a ground-truth trajectory in which earlier tasks provide essential information or strategies for later ones. At each step t , the agent processes input x_t , retrieves and synthesizes memory, produces prediction \hat{y}_t , and updates the memory state M_t , yielding the predicted trajectory:

$$(x_1, \hat{y}_1, M_1) \rightarrow (x_2, \hat{y}_2, M_2) \rightarrow \dots \rightarrow (x_T, \hat{y}_T, M_T).$$

This design transforms static benchmarks into interactive evaluation streams that explicitly probe an LLM’s ability to accumulate, adapt, and refine knowledge during deployment.

3.2. ExpRAG: Experience Retrieval and Aggregation

As a simple baseline and extension, we define **ExpRAG**, a task-level retrieval-augmented agent. Each memory entry $m_i = S(x_i, \hat{y}_i, f_i)$ encodes a structured experience text with template S . At step t , the agent retrieves k similar experiences from memory according to a retrieval score ϕ :

$$R_t = \text{Top-}k_{m_i \in M_t} \phi(x_t, m_i).$$

The model conditions on these retrieved examples following the in-context learning principle:

$$\hat{y}_t = F(x_t, R_t),$$

and appends the new experience to memory:

$$M_{t+1} = M_t \cup \{(x_t, \hat{y}_t, f_t)\}.$$

ExpRAG thus performs one-shot experience reuse through retrieval and aggregation. It captures how simple memory-based extensions of in-context learning behave but lacks iterative reasoning or adaptive refinement during inference.

3.3. ReMem: Synergizing Reasoning, Acting, and Memory

We propose **ReMem**, a simple yet effective framework that unifies reasoning, action, and memory refinement within a single decision loop. Unlike conventional retrieval-augmented or ReAct-style methods that treat memory as static context, ReMem introduces a third dimension of *memory reasoning*, allowing the agent to actively evaluate, reorganize, and evolve its own memory during problem solving.

At each step t , given the current input x_t , memory state M_t , and previous reasoning traces $o_t^{1:n-1}$ at this step, the agent selects one of three operations:

$$a_t^n \in \{\text{Think, Act, Refine}\}.$$

It then performs the operation and transitions according to:

$$o_t^n = \text{Agent}(x_t, M_t, a_t^n),$$

where o_t^n denotes the output generated at step t after n operations, such as an intermediate reasoning trace, an external action, and memory refine thoughts.

Specifically, *Think* produces internal reasoning traces that help decompose the task and guide subsequent actions; *Act* executes an operation in the environment or outputs a response observable to the user; *Refine* performs meta-reasoning over memory, which exploiting useful experiences, pruning noise, and reorganizing M_t , to better support future reasoning and action. Within each step, the agent may perform multiple rounds of *Think* and *Refine*, and the step terminates once an *Act* operation is selected. This induces a Markov decision process where the state at step t after n operations is $s_t^n = (x_t, M_t, o_t^{1:n-1})$, the action space is $\{\text{Think, Act, Refine}\}$, and the transition dynamics are given by the Agent operator together with the environment response. Depending on the task, the *Act*-output of step t may serve as the final answer for single-step tasks or as an intermediate result in multi-step settings, where the process continues until the overall task is completed.

This unified formulation expands the action space of ReAct-style (Yao et al., 2022) agents by introducing an explicit memory reasoning mechanism. Through this extension, memory becomes an adaptive component that interacts with reasoning in real time rather than remaining a passive context. Under this view, the entire decision loop can also be interpreted as a Markov process, where the state

encapsulates the current input, memory state, and ongoing reasoning traces. Such integration yields a lightweight yet powerful paradigm for continual adaptation, where the agent learns to reason about both the task and its own knowledge state. By coupling reflection with memory evolution, ReMem establishes a new standard for adaptive, self-improving LLM agents.

4. Experiments

In this section, we evaluate leading LLMs on the Evo-Memory benchmark under our unified test-time learning pipeline, focusing on five key research questions (RQs):

- **RQ1:** How do LLM agents perform on Evo-Memory across diverse domains and task types, and does REMEM enhance their test-time learning ability?
- **RQ2:** What factors influence the effectiveness of memory in different tasks, and how does experience reuse improve task efficiency?
- **RQ3:** How does task sequence difficulty (e.g., easy vs. hard trajectories) affect memory adaptation and generalization?
- **RQ4:** How do varying feedback types impact learning dynamics and memory refinement across tasks?
- **RQ5:** How does cumulative performance evolve over task sequences and time steps, reflecting continual adaptation during deployment?

4.1. Experimental Setup

Evo-Memory evaluates memory mechanisms under realistic streaming multi-task conditions. In what follows, we describe the benchmark datasets, metrics, and the methods compared.

4.1.1. Datasets

Evo-Memory is evaluated on a diverse suite of datasets spanning factual knowledge, reasoning, mathematics, programming, and goal-oriented interaction. For factual and reasoning ability, we include **MMLU-Pro** (Zheng et al., 2024) and **GPQA-Diamond** (Rein et al., 2024), which test multi-disciplinary and graduate-level reasoning. For mathematical problem solving, we use **AIME-24** and **AIME-25** (HuggingFaceH4, 2024), containing Olympiad-style challenges requiring symbolic reasoning and exact-match evaluation. For tool-use and API grounding, we include **ToolBench** (Patil et al., 2023). For multi-turn and goal-oriented interaction, we adopt the **AgentBoard** (Zhuang et al., 2024) suite, covering **AlfWorld** (Shridhar et al., 2021), **BabyAI** (Chevalier-Boisvert et al., 2019), **ScienceWorld** (Wang et al., 2022), **Jericho** (Hausknecht et al., 2020), and **PDDL** tasks (Yang et al., 2023). Together, these datasets span both single-turn and interactive settings, enabling a unified evaluation of factual recall, procedural reasoning, and long-horizon adaptation. All methods are evaluated under the same *search–predict–evolve* loop

$$(x_t, M_t) \xrightarrow{\text{search}} R_t \xrightarrow{\text{synthesis}} \hat{y}_t \xrightarrow{\text{evolve}} M_{t+1},$$

with identical prompting templates, configurations, and memory budgets if applicable. Feedback f_t is considered as the correctness signal.

4.1.2. Evaluation

Evo-Memory evaluates both task performance and memory quality along four dimensions. First, **answer accuracy** measures whether the model produces correct outputs in single-turn tasks. Second,

success rate and **progress rate** evaluate goal completion in multi-turn tasks. Third, **step efficiency** tracks how many steps are needed to reach a goal, reflecting reasoning conciseness. Finally, **sequence robustness** tests whether performance stays stable under different task orders. Together, these metrics assess how well the agent learns, adapts, and reuses knowledge over time.

4.1.3. Methods

We benchmark a broad range of agents and memory architectures instantiated on two strong **LLM backbones**: the Gemini-2.5 series (Comanici et al., 2025) (FLASH, FLASH-LITE, and PRO) and the Claude family (Anthropic, 2025) (3.5-HAIKU and 3.7-SONNET). The evaluated methods are grouped into four categories: (1) **Agent pipelines without persistent memory**, including ReAct (Yao et al., 2022) and Amem (Xu et al., 2025), which rely on short-term context or lightweight caches; (2) **Adaptive agentic memory methods**, such as SelfRAG (Asai et al., 2024b), MemOS (Li et al., 2025), Mem0 (Chhikara et al., 2025), and LangMem (LangChain contributors, 2025), which support dynamic retrieval and continual updates; (3) **Memory-based agents for procedural knowledge**, including Dynamic Cheatsheet (DC) (Suzgun et al., 2025) with two variants Cumulative (Cu) and Synthesis (RS) and Agent Workflow Memory (AWM) (Wang et al., 2024), which emphasize reusable workflows and task strategies; and (4) **Proposed evolving-memory framework**, comprising ExpRecent, ExpRAG, and ReMem, which unify reasoning, action, and memory refinement in a self-evolving loop. All methods are evaluated under a unified *search–predict–evolve* protocol to isolate the effects of memory design. Implementation and prompting details are provided in Appendix A. We exclude systems such as MemoryGpt (Zhong et al., 2023) and MemoryBank (Zhong et al., 2023) that target factual recall only, since Evo-Memory is designed to test evolving and procedural memory. Our goal is not to improve or modify the underlying LLMs themselves. Certain methods, such as MemOS and LangMem, are not fully compatible with embodied environments, and thus we exclude them from multi-turn datasets. Evo-Memory isolates the effect of search and evolution in *memory mechanisms*, so that observed differences reflect solely memory design rather than raw LLM capability.

4.2. Experiments

Below are the conducted experiments to answer the proposed research questions.

4.3. Analysis of Results (RQ1)

Tables 1 and 2 summarize the results across single-turn and multi-turn settings. Overall, Evo-Memory demonstrates that self-evolving memory architectures provide consistent improvements. In single-turn reasoning and QA benchmarks (AIME-24/25, GPQA, MMLU-Pro, ToolBench), evolving-memory methods show consistent improvements, with ReMem achieving 0.65 average exact match and 0.85/0.71 API accuracy under Gemini-2.5 Flash. Adaptive retrieval methods enhance factual grounding, yet only evolving systems maintain consistent gains through iterative refinement. Agents with procedural knowledge perform well on structured domains such as AIME but lag in scientific reasoning and tool use, showing limited flexibility. ExpRAG serves as a simple yet highly effective baseline, outperforming several more complex designs. While improvements in single-turn settings are moderate, the overall trend remains consistent across datasets and model families.

In multi-turn reasoning environments (AlfWorld, BabyAI, PDDL, ScienceWorld), ReMem and ExpRAG achieve strong and stable performance on both Gemini-2.5 and Claude backbones, reaching 0.92/0.96 on BabyAI and 0.95/0.62 on ScienceWorld. These results indicate that continual reflection and refinement substantially improve procedural knowledge accumulation. Performance gains are notably larger in multi-turn settings, underscoring that continual adaptation becomes

LLM Backbone	Method	Exact Match \uparrow						API / Acc. \uparrow	
		AIME24	AIME25	GPQA	MMLU-Pro (Eco.)	MMLU-Pro (Eng.)	MMLU-Pro (Philo.)	ToolBench	Avg. \uparrow
Claude 3.7 Sonnet	Baseline	0.17	0.13	0.55	0.84	0.63	0.78	0.76/0.62	0.54
	History	0.13	0.23	0.56	0.85	0.64	0.78	0.76/0.61	0.55
	ReAct	0.17	0.10	0.57	0.84	0.63	0.76	0.76/0.61	0.54
	Amem	0.27	0.17	0.54	0.83	0.63	0.79	0.77/0.63	0.56
	SelfRAG	0.20	0.10	0.58	0.84	0.65	0.77	0.77/0.63	0.55
	MemOS	0.17	0.20	0.55	0.84	0.64	0.76	0.76/0.62	0.55
	Mem0	0.20	0.13	0.58	0.84	0.62	0.77	0.76/0.61	0.55
	LangMem	0.10	0.13	0.53	0.77	0.56	0.66	0.77/0.63	0.49
	DC-Cu	0.17	0.23	0.57	0.79	0.52	0.65	0.77/0.62	0.52
	DC-RS	0.20	0.20	0.62	0.79	0.52	0.60	0.77/0.62	0.52
	AWM	0.03	0.03	0.53	0.80	0.56	0.72	0.76/0.62	0.48
	ExpRecent	0.13	0.20	0.61	0.86	0.63	0.78	0.82/0.66	0.56
	ExpRAG	0.17	0.17	0.70	0.85	0.67	0.80	0.88/0.72	0.59
	ReMem	0.13	0.13	0.67	0.86	0.65	0.80	0.87/0.71	0.58
	Baseline	0.47	0.47	0.48	0.83	0.46	0.75	0.71/0.61	0.59
	History	0.60	0.47	0.43	0.84	0.42	0.78	0.62/0.54	0.58
Gemini 2.5 Flash	ReAct	0.30	0.27	0.05	0.64	0.16	0.54	0.64/0.57	0.37
	Amem	0.70	0.57	0.52	0.83	0.42	0.72	0.72/0.60	0.63
	SelfRAG	0.50	0.47	0.46	0.83	0.45	0.75	0.72/0.61	0.59
	MemOS	0.47	0.47	0.50	0.82	0.46	0.75	0.71/0.61	0.59
	Mem0	0.50	0.47	0.45	0.83	0.46	0.74	0.71/0.61	0.59
	LangMem	0.43	0.50	0.53	0.79	0.39	0.71	0.68/0.57	0.57
	DC-Cu	0.60	0.40	0.48	0.79	0.44	0.69	0.70/0.59	0.58
	DC-RS	0.53	0.37	0.48	0.80	0.42	0.69	0.68/0.57	0.56
	AWM	0.50	0.37	0.49	0.79	0.43	0.72	0.71/0.59	0.56
	ExpRecent	0.47	0.47	0.42	0.83	0.39	0.75	0.78/0.66	0.58
	ExpRAG	0.43	0.47	0.42	0.83	0.43	0.78	0.87/0.73	0.60
	ReMem	0.60	0.53	0.51	0.85	0.46	0.79	0.85/0.71	0.65

Table 1 | Cross-dataset results of diverse memory architectures across models on the single-turn reasoning and question answering datasets. Categories are separated by horizontal rules; results (Exact Match \uparrow and API/Acc \uparrow) compare zero-shot, agentic, adaptive, procedural, and proposed memory methods.

increasingly valuable as task horizons lengthen. While many baselines enhance retrieval grounding, they struggle to reuse long-horizon experiences and often falter in open-ended environments. Notably, lightweight variants such as `ExpRecent` and `ExpRAG` still perform competitively despite their simplicity, suggesting that explicit task-level utilization during test-time evolution is both promising and underexplored.

Across all experiments, evolving-memory methods demonstrate consistent gains on both Gemini and Claude backbones. Smaller models benefit particularly from self-evolving memory, suggesting that test-time refinement is a practical path to enhancing the capability of lighter LLMs. Together, these findings establish task-level memory utilization and continual reorganization as valuable directions for future research, providing a standardized reference point for developing and evaluating evolving-memory agents. Additional results across more LLM families are presented in Appendix B.1.

4.4. Analysis of Memory Improvement (RQ2)

Figure 4 shows that `ReMem`’s improvement strongly correlates with within-dataset task similarity (Pearson $r = 0.717$ on Gemini 2.5 Flash and $r = 0.563$ on Claude 3.7 Sonnet). Task similarity is measured by computing the average cosine distance between each task embedding and its dataset cluster center, where embeddings are obtained from the retriever encoder. A smaller average distance

LLM Backbone	Method	Alf World		BabyAI		PDDL		ScienceWorld		Avg.	
		S	P	S	P	S	P	S	P	S	P
Gemini 2.5 Flash	Baseline	0.12	0.34	0.61	0.71	0.12	0.20	0.24	0.59	0.27	0.46
	History	0.28	0.60	0.52	0.64	0.08	0.15	0.31	0.71	0.30	0.53
	ReAct	0.24	0.56	0.48	0.63	0.22	0.33	0.34	0.71	0.32	0.56
	Amem	0.25	0.59	0.53	0.64	0.10	0.16	0.36	0.74	0.31	0.53
	SelfRAG	0.25	0.59	0.52	0.65	0.08	0.16	0.34	0.74	0.30	0.54
	Mem0	0.27	0.61	0.54	0.66	0.10	0.19	0.32	0.70	0.31	0.54
	DC-Cu	0.25	0.59	0.53	0.64	0.08	0.17	0.29	0.71	0.29	0.53
	DC-RS	0.27	0.60	0.53	0.66	0.07	0.15	0.33	0.73	0.30	0.54
	AWM	0.26	0.59	0.52	0.64	0.08	0.16	0.33	0.73	0.30	0.53
	ExpRecent	0.37	0.65	0.53	0.64	0.13	0.22	0.53	0.83	0.39	0.59
	ExpRAG	0.59	0.79	0.56	0.65	0.17	0.27	0.53	0.81	0.46	0.63
	ReMem	0.66	0.81	0.53	0.61	0.22	0.33	0.58	0.81	0.50	0.64
Claude 3.7 Sonnet	Baseline	0.18	0.49	0.51	0.66	0.17	0.39	0.10	0.53	0.24	0.52
	History	0.50	0.73	0.48	0.66	0.65	0.85	0.32	0.74	0.49	0.74
	ReAct	0.51	0.75	0.57	0.72	0.75	0.91	0.44	0.77	0.57	0.79
	Amem	0.48	0.73	0.46	0.64	0.62	0.84	0.33	0.73	0.47	0.73
	SelfRAG	0.52	0.75	0.46	0.64	0.65	0.84	0.31	0.74	0.49	0.74
	Mem0	0.51	0.74	0.48	0.66	0.65	0.84	0.37	0.76	0.50	0.75
	DC-Cu	0.50	0.74	0.50	0.67	0.62	0.84	0.33	0.75	0.49	0.75
	DC-RS	0.50	0.74	0.52	0.68	0.62	0.84	0.34	0.74	0.50	0.75
	AWM	0.49	0.73	0.53	0.68	0.60	0.82	0.34	0.74	0.49	0.74
	ExpRecent	0.66	0.83	0.63	0.73	0.53	0.76	0.49	0.82	0.58	0.79
	ExpRAG	0.74	0.89	0.62	0.72	0.72	0.89	0.46	0.76	0.63	0.82
	ReMem	0.92	0.96	0.73	0.83	0.83	0.95	0.62	0.89	0.78	0.91

Table 2 | Cross-environment results across four embodied reasoning benchmarks (Alf World, BabyAI, PDDL, ScienceWorld). Each dataset reports success (S) and progress (P) rates. Bold indicates the best (including ties) per column. The last two columns show averaged S and P across datasets.

indicates higher intra-dataset coherence and thus stronger structural similarity. Tasks with higher embedding cluster ratios, such as PDDL and Alf World, yield larger gains, suggesting that recurring task structures facilitate memory reuse and generalization. In contrast, more diverse or low-similarity datasets like AIME-25 or GPQA show smaller gains, reflecting limited transferable experiences. These findings highlight the importance of embedding organization and semantic overlap in driving effective memory evolution. Further analysis of memory pruning rates can be found in Appendix B.2.

Figure 5 compares step efficiency across four environments. Evolving-memory methods consistently require fewer steps to reach completion, with REMEM achieving the strongest and most stable reductions (e.g., from 22.6 to 11.5 steps on Alf World). The lightweight EXP RAG and EXP RECENT also perform competitively, showing that simple task-level evolution can greatly improve efficiency without extra complexity. Overall, continual refinement not only boosts accuracy but also makes reasoning more focused and efficient.

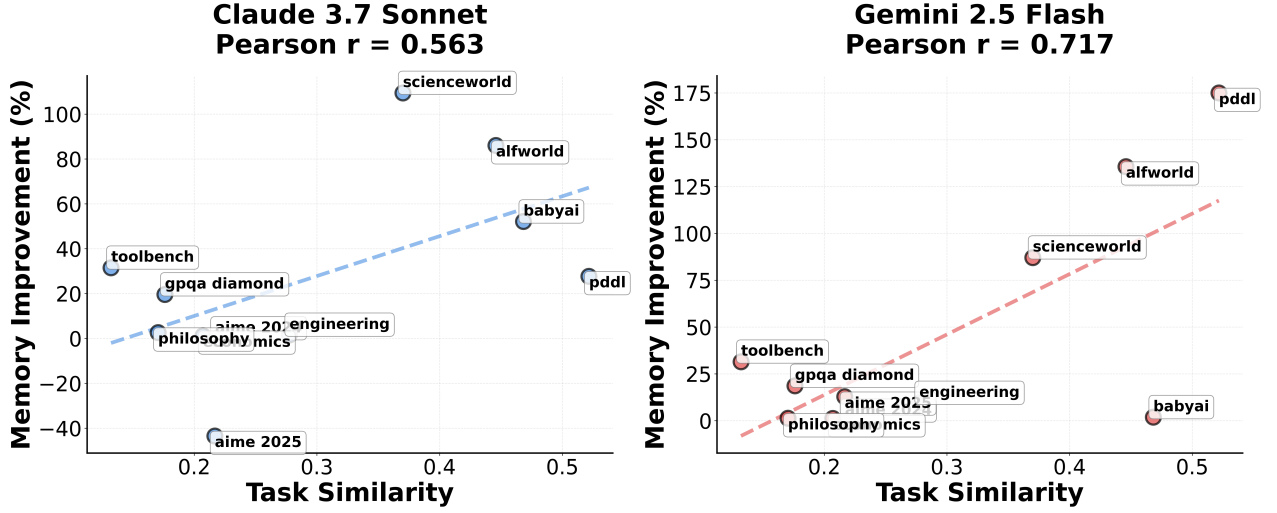


Figure 4 | ReMem performance gain over history baseline versus within-dataset task similarity.

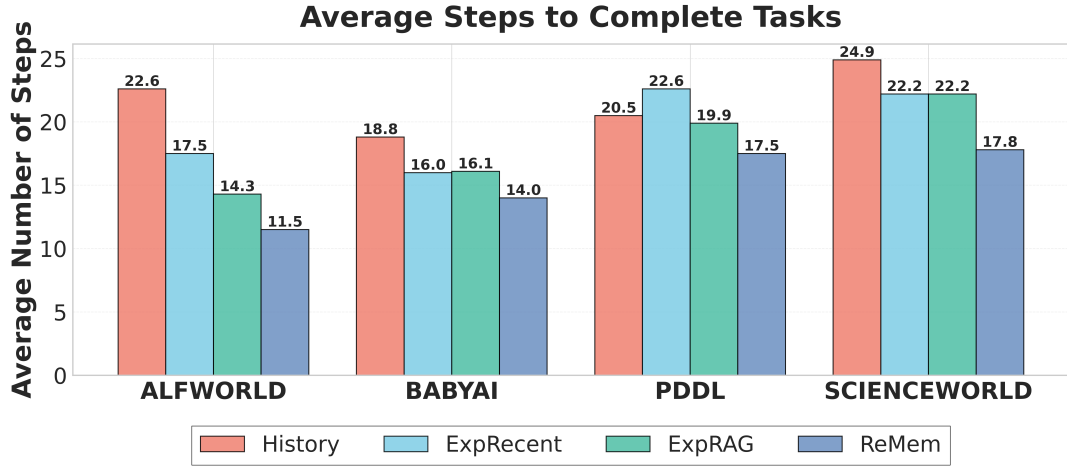


Figure 5 | Average steps to complete tasks across four benchmarks. We compare four methods: History, ExpRecent, ExpRAG, and ReMem. Lower is better. ReMem consistently requires fewer steps to complete tasks across all datasets, demonstrating more efficient task execution..

Direction	Method	Alf World		ScienceWorld		Avg.	
		S	P	S	P	S	P
Easy→Hard	Base	0.50	0.73	0.32	0.74	0.41	0.74
	ExpRecent	0.66	0.82	0.48	0.83	0.57	0.83
	ExpRAG	0.77	0.87	0.37	0.71	0.57	0.79
	ReMem	0.91	0.96	0.63	0.88	0.77	0.92
Hard→Easy	ExpRecent	0.72	0.85	0.47	0.80	0.60	0.83
	ExpRAG	0.87	0.92	0.51	0.81	0.69	0.87
	ReMem	0.94	0.97	0.68	0.90	0.81	0.94

Table 3 | Comparison of memory-based agents under different sequence difficulty directions. Each cell reports Success (S) and Progress (P). Easy→Hard and Hard→Easy indicate task order transitions, and averages (Avg) summarize per-direction performance.

Model	Method	Alf World		ScienceWorld		Avg.	
		S	P	S	P	S	P
Claude 3.7 Sonnet	Amem	0.49	0.73	0.31	0.74	0.40	0.74
	SelfRAG	0.47	0.73	0.34	0.73	0.41	0.73
	Mem0	0.49	0.73	0.36	0.74	0.43	0.74
	DC-Cu	0.52	0.75	0.34	0.73	0.43	0.74
	DC-RS	0.51	0.74	0.38	0.74	0.45	0.74
	AWM	0.55	0.76	0.32	0.72	0.44	0.74
	ExpRecent	0.62	0.80	0.34	0.74	0.48	0.77
	ExpRAG	0.76	0.90	0.27	0.63	0.52	0.77
	ReMem	0.92	0.96	0.69	0.91	0.81	0.94
Gemini 2.5 Flash	Amem	0.22	0.57	0.39	0.75	0.31	0.66
	SelfRAG	0.25	0.58	0.36	0.71	0.31	0.65
	Mem0	0.25	0.59	0.34	0.71	0.30	0.65
	DC-Cu	0.20	0.56	0.36	0.72	0.28	0.64
	DC-RS	0.21	0.57	0.36	0.71	0.29	0.64
	AWM	0.19	0.56	0.36	0.74	0.28	0.65
	ExpRecent	0.22	0.57	0.59	0.86	0.41	0.72
	ExpRAG	0.25	0.60	0.51	0.78	0.38	0.69
	ReMem	0.57	0.76	0.50	0.75	0.54	0.76

Table 4 | Results with both successful and failed task experiences on Alf World and ScienceWorld. Each cell reports Success (S) and Progress (P). Horizontal rules separate method families. Bold numbers denote the best results per metric within each model.

4.5. Task Sequence: Easy v.s. Hard (RQ3)

Table 3 examines how memory-based agents adapt to changes in task difficulty. Baseline methods exhibit noticeable degradation when moving from easier to harder tasks, revealing limited robustness under distribution shifts. In contrast, evolving-memory agents, particularly ReMem, maintain strong and consistent performance across both directions, reaching up to 0.94/0.97 success and progress in the Hard→Easy setting. This stability shows that continual reflection enables ReMem to retain transferable knowledge even as task complexity varies. The results further indicate that the design of task sequences, especially the ordering of difficulty levels, has a substantial impact on evaluating memory adaptation. They also suggest that well-structured task progressions can actively facilitate learning by allowing models to build on prior experiences and generalize across increasingly complex challenges. Together, these findings highlight the importance of standardized and thoughtfully organized task sequences for both fair evaluation and effective model development in future benchmarks.

4.6. Analysis of Feedback (RQ4)

Table 4 evaluates how agents perform when both successful and failed task experiences are stored in memory. Baseline methods experience a clear performance drop when exposed to unfiltered failures, indicating that naive memory accumulation introduces noise and disrupts subsequent retrieval. In contrast, evolving-memory approaches, particularly ReMem, remain robust by actively refining stored experiences, achieving the highest overall success and progress rates under both Claude and Gemini backbones. These results demonstrate that selective utilization, which involves learning from successes while appropriately leveraging information from failures, is crucial for stable test-time adaptation. They further highlight that memory refinement plays a central role in handling imperfect experiences and suggest that future work should explore failure-aware strategies for memory evolution.

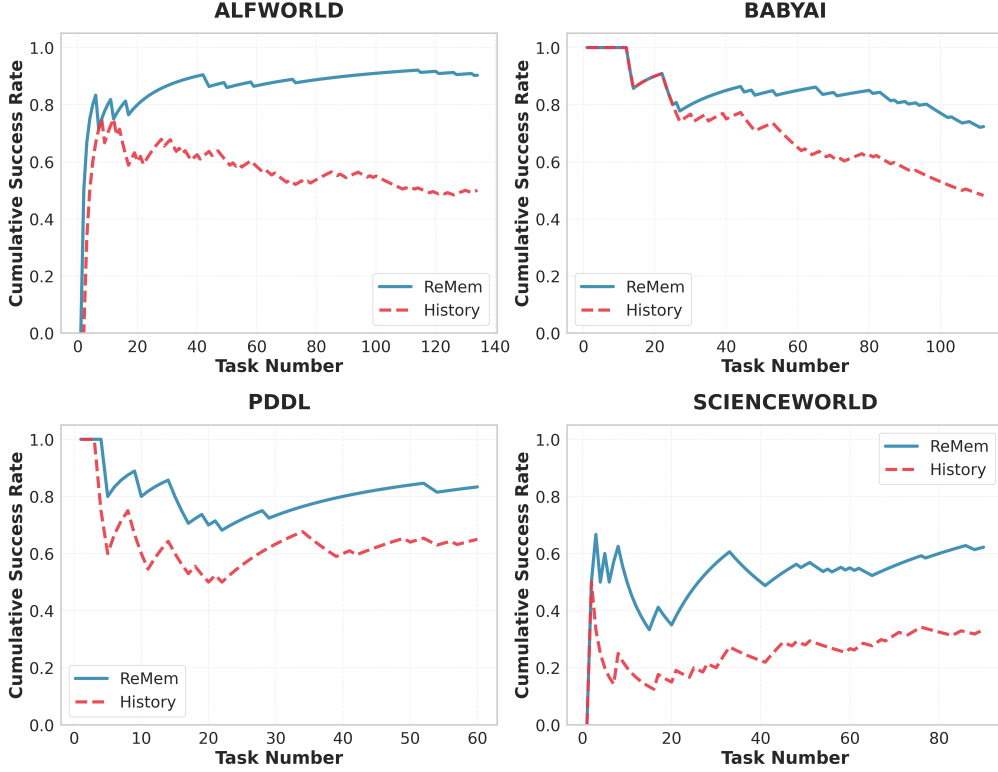


Figure 6 | Cumulative success rate across four interactive agent datasets. ReMem (solid blue) outperforms the History baseline (dashed red) on ALFWorld, BabyAI, PDDL, and ScienceWorld tasks. The rolling average shows performance trends as more task instances are evaluated.

4.7. Performance w.r.t Time Steps (RQ5)

Figure 6 shows the cumulative accuracy as tasks progress across four interactive environments. The curves mainly serve to compare REMEM with the History baseline, as individual trajectories carry limited standalone meaning. Across all environments, REMEM consistently achieves faster adaptation and more stable retention over time. These results highlight that continual reflection enables REMEM to sustain performance across long task sequences, illustrating its robustness in test-time learning. More comparative results on single-turn tasks are presented in Appendix B.3.

5. Conclusion

Self-evolving memory is a fundamental yet underexplored aspect of LLM capability. While prior work centers on static conversational recall, it overlooks how models accumulate and reuse experience across evolving task streams. Evo-Memory fills this gap by transforming static datasets into streaming trajectories, systematically evaluating how LLMs retrieve, adapt, and refine memory through interaction. Our results show that memory can substantially enhance performance but remains fragile in stability and procedural reuse. To foster progress, we introduce ExpRAG for experience retrieval and ReMem for interleaving reasoning, action, and memory updates. We hope Evo-Memory serves as a unified platform for building LLMs with reliable and continually improving memory.

References

- Anthropic. Introducing Claude 4: Claude Opus 4 and Claude Sonnet 4. <https://www.anthropic.com/news/claude-4>, 2025. Accessed: 2025-09-10.
- A. Asai, Z. Wu, Y. Wang, A. Sil, and H. Hajishirzi. Self-rag: Self-reflective retrieval augmented generation. In *NeurIPS 2023 workshop on instruction tuning and instruction following*, 2023.
- A. Asai, Z. Wu, Y. Wang, A. Sil, and H. Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024a. URL <https://openreview.net/forum?id=hSyW5go0v8>.
- A. Asai, Z. Wu, Y. Wang, A. Sil, and H. Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. 2024b.
- J. Chen, S. Xiao, P. Zhang, K. Luo, D. Lian, and Z. Liu. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation, 2023.
- R. Chen, Z. Li, and Y. Wang. Llm-as-optimizer: Self-improving agents through differentiable feedback. *arXiv preprint arXiv:2503.04567*, 2025.
- M. Chevalier-Boisvert, D. Bahdanau, S. E.-T. Lahlou, L. Willems, H. Lozano, L. Dassa, S. Kim, J. Pineau, and A. Courville. Babyai: A platform to study the sample efficiency of grounded language learning. In *International Conference on Learning Representations (ICLR)*, 2019.
- P. Chhikara, D. Khant, S. Aryan, T. Singh, and D. Yadav. Mem0: Building production-ready ai agents with scalable long-term memory. *arXiv preprint arXiv:2504.19413*, 2025.
- G. Comanici, E. Bieber, M. Schaekermann, I. Pasupat, N. Sachdeva, I. Dhillon, M. Blistein, O. Ram, D. Zhang, E. Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Y. Gao, Z. Sun, J. Huang, H. Zhang, Y. Wang, Y. Luo, and C. Finn. A survey of self-evolving agents: On path to artificial super intelligence. *arXiv preprint arXiv:2507.21046*, 2025.
- M. Hausknecht, P. Ammanabrolu, M.-A. Côté, and X. Yuan. Interactive fiction games: A colossal adventure for reinforcement learning agents. In *AAAI Conference on Artificial Intelligence*, 2020.
- Y. Hu, Y. Wang, and J. McAuley. Evaluating memory in llm agents via incremental multi-turn interactions. *arXiv preprint arXiv:2507.05257*, 2025.
- J. Huang, R. Zhao, Y. Li, and Y. Zhou. Self-discovering agents: Autonomous skill expansion via continual interaction. *arXiv preprint arXiv:2506.08791*, 2025.
- HuggingFaceH4. AIME-24: American Invitational Mathematics Examination 2024 Benchmark. https://huggingface.co/datasets/HuggingFaceH4/aime_2024, 2024. Accessed: 2025-09-10.
- HuggingFaceH4. AIME-25: American Invitational Mathematics Examination 2025 Benchmark. https://huggingface.co/datasets/HuggingFaceH4/aime_2025, 2025. Accessed: 2025-09-10.
- Y. Iwasawa and Y. Matsuo. T3a: Test-time template adjustments for domain generalization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

- LangChain contributors. Langchain, 2025. URL <https://github.com/langchain-ai/langchain>. MIT License.
- P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. Yih, T. Rocktäschel, S. Riedel, and D. Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>.
- Z. Li, S. Song, H. Wang, S. Niu, D. Chen, J. Yang, C. Xi, H. Lai, J. Zhao, Y. Wang, J. Ren, Z. Lin, J. Huo, T. Chen, K. Chen, K.-R. Li, Z. Yin, Q. Yu, B. Tang, H. Yang, Z. Xu, and F. Xiong. Memos: An operating system for memory-augmented generation (mag) in large language models. *ArXiv*, abs/2505.22101, 2025. URL <https://api.semanticscholar.org/CorpusID:278960153>.
- X. Liang, B. Wang, H. Huang, S. Wu, P. Wu, L. Lu, Z. Ma, and Z. Li. Scm: Enhancing large language model with self-controlled memory framework. 2023. URL <https://api.semanticscholar.org/CorpusID:258331553>.
- J. Liu. LlamaIndex, 11 2022. URL https://github.com/jerryjliu/llama_index.
- J. Liu, N. Loo, H. Li, R. Chen, X. Chen, T. Hospedales, and Y. Wang. Ttt++: When does test-time training fail or thrive? In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.
- A. Maharana, D.-H. Lee, S. Tulyakov, M. Bansal, F. Barbieri, and Y. Fang. Evaluating very long-term conversational memory of llm agents. *ArXiv*, abs/2402.17753, 2024a. URL <https://api.semanticscholar.org/CorpusID:268041615>.
- A. Maharana, D.-H. Lee, S. Tulyakov, M. Bansal, F. Barbieri, and Y. Fang. Evaluating very long-term conversational memory of llm agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13851–13870, 2024b.
- A. Modarressi, A. Imani, M. Fayyaz, and H. Schütze. Ret-llm: Towards a general read-write memory for large language models. *ArXiv*, abs/2305.14322, 2023. URL <https://api.semanticscholar.org/CorpusID:258841042>.
- Y. Niu, Z. Li, B. Du, and D. Tao. Efficient test-time adaptation via sample-efficient entropy minimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- S. Ouyang, W. Yu, K. Ma, Z.-Q. Xiao, Z. Zhang, M. Jia, J. Han, H. Zhang, and D. Yu. Repograph: Enhancing ai software engineering with repository-level code graph. *ArXiv*, abs/2410.14684, 2024. URL <https://api.semanticscholar.org/CorpusID:273502041>.
- C. Packer, V. Fang, S. G. Patil, K. Lin, S. Wooders, and J. Gonzalez. Memgpt: Towards llms as operating systems. *ArXiv*, abs/2310.08560, 2023. URL <https://api.semanticscholar.org/CorpusID:263909014>.
- J. S. Park, C. O’Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein. Generative agents: Interactive simulacra of human behavior. In *ACM Symposium on User Interface Software and Technology (UIST)*, 2023.
- S. G. Patil, H. Li, T. Zhang, et al. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023.

- P. Rasmussen, P. Paliychuk, T. Beauvais, J. Ryan, and D. Chalef. Zep: A temporal knowledge graph architecture for agent memory. *ArXiv*, abs/2501.13956, 2025. URL <https://api.semanticscholar.org/CorpusID:275907122>.
- D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
- N. Shinn, F. Cassano, B. Labash, A. Gopinath, K. Narasimhan, and S. Yao. Reflexion: language agents with verbal reinforcement learning. In *Neural Information Processing Systems*, 2023. URL <https://api.semanticscholar.org/CorpusID:258833055>.
- M. Shridhar, J. Thomason, D. Gordon, Y. Bisk, W. Han, R. Mottaghi, L. Zettlemoyer, and D. Fox. Alfworld: Aligning text and embodied environments for interactive learning. In *International Conference on Learning Representations (ICLR)*, 2021.
- M. Suzgun, M. Yuksekgonul, F. Bianchi, D. Jurafsky, and J. Zou. Dynamic cheatsheet: Test-time learning with adaptive memory. *arXiv preprint arXiv:2504.07952*, 2025.
- D. Wang, E. Shelhamer, S. Liu, B. Olshausen, and T. Darrell. Tent: Fully test-time entropy minimization. In *International Conference on Learning Representations (ICLR)*, 2021.
- G. Wang, Y. Wang, Z. Wu, G. Chen, Z. Huang, H. Zhao, S. Han, V. Koltun, J. Zhu, and K. Lin. Voyager: An open-ended embodied agent with large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- Y. Wang, L. Yuan, K. Gopalakrishnan, A. Narayan-Chen, A. Fang, and M. Hausknecht. Scienceworld: Is your agent smarter than a 5th grader? In *NeurIPS*, 2022.
- Z. Z. Wang, J. Mao, D. Fried, and G. Neubig. Agent workflow memory. *ArXiv*, abs/2409.07429, 2024. URL <https://api.semanticscholar.org/CorpusID:272592995>.
- C.-K. Wu, Z. R. Tam, C.-Y. Lin, Y.-N. V. Chen, and H.-y. Lee. Streambench: Towards benchmarking continuous improvement of language agents. *Advances in Neural Information Processing Systems*, 37:107039–107063, 2024a.
- D. Wu, H. Wang, W. Yu, Y. Zhang, K.-W. Chang, and D. Yu. Longmemeval: Benchmarking chat assistants on long-term interactive memory. In *The Thirteenth International Conference on Learning Representations*.
- D. Wu, H. Wang, W. Yu, Y. Zhang, K.-W. Chang, and D. Yu. Longmemeval: Benchmarking chat assistants on long-term interactive memory. *ArXiv*, abs/2410.10813, 2024b. URL <https://api.semanticscholar.org/CorpusID:273345961>.
- W. Xu, K. Mei, H. Gao, J. Tan, Z. Liang, and Y. Zhang. A-mem: Agentic memory for llm agents. *arXiv preprint arXiv:2502.12110*, 2025.
- S. Yan, X. Yang, Z. Huang, E. Nie, Z. Ding, Z. Li, X. Ma, H. Schütze, V. Tresp, and Y. Ma. Memory-r1: Enhancing large language model agents to manage and utilize memories via reinforcement learning. *arXiv preprint arXiv:2508.19828*, 2025.
- R. Yang, Q. Zhou, Y. Zhang, J. Zhang, X. Zhang, K. Zhang, S. Xu, W. Chen, H. Ma, W. Wang, et al. Pddlbench: Benchmarking llms on symbolic planning with pddl. *arXiv preprint arXiv:2312.00754*, 2023.

- S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. R. Narasimhan, and Y. Cao. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*, 2022.
- S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- H. Yu, T. Chen, J. Feng, J. Chen, W. Dai, Q. Yu, Y.-Q. Zhang, W.-Y. Ma, J. Liu, M. Wang, et al. Memagent: Reshaping long-context llm with multi-conv rl-based memory agent. *arXiv preprint arXiv:2507.02259*, 2025.
- M. Zhang, K. Ahuja, K.-C. Lee, T. Zhang, and C. Finn. Memo: Test-time robustness via adaptation and augmentation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- R. Zhao, Y. Li, Z. Qian, X. Wang, W. Zhao, and J. Huang. Self-evolving agents: Continual test-time learning through memory and reflection. *arXiv preprint arXiv:2507.21046*, 2025.
- J. Zheng, X. Cai, Q. Li, D. Zhang, Z. Li, Y. Zhang, L. Song, and Q. Ma. Lifelongagentbench: Evaluating llm agents as lifelong learners. *arXiv preprint arXiv:2505.11942*, 2025.
- Y. Zheng, T. Li, H. Li, C. Zhao, J. Wang, Z. Zhang, S. Deng, N. Zhang, and H. Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024.
- W. Zhong, L. Guo, Q.-F. Gao, H. Ye, and Y. Wang. Memorybank: Enhancing large language models with long-term memory. *ArXiv*, abs/2305.10250, 2023. URL <https://api.semanticscholar.org/CorpusID:258741194>.
- Y. Zhou, Z. Sun, J. Huang, K.-H. Lee, Y. Luo, and C. Finn. Eureka: Human-level reward design via coding large language models. *arXiv preprint arXiv:2402.10654*, 2024.
- Z. Zhou, A. Qu, Z. Wu, S. Kim, A. Prakash, D. Rus, J. Zhao, B. K. H. Low, and P. P. Liang. MEM1: Learning to synergize memory and reasoning for efficient long-horizon agents. *arXiv preprint arXiv:2506.15841*, 2025.
- H. Zhuang, B. Zhang, Y. Wang, Y. Xie, Z. Wang, Q. Zhang, H. Zhang, C. Zhang, X. Zhou, J. He, et al. Agentboard: Evaluating long-term memory and generalization in large language model agents. *arXiv preprint arXiv:2401.13178*, 2024.

Appendix

Contents

A	Experimental Details	19
A.1	Datasets	19
A.2	Configuration	19
A.3	Evaluation	20
A.4	Methods	20
B	Experiments	21
B.1	Additional Experiments	21
B.2	Additional Analysis of Memory Pruning	22
B.3	Additional Comparative Curves on Single-turn Tasks	22
C	Prompts	23
D	Limitations	24
E	Use of Large Language Models	24

A. Experimental Details

Evo-Memory evaluates memory mechanisms under realistic streaming multi-task conditions. In what follows, we describe the benchmark datasets, metrics, configurations and the methods compared in details.

A.1. Datasets

We evaluate our approach on a diverse suite of benchmarks that span factual knowledge, reasoning, mathematics, programming, and goal-oriented interaction.

We first introduce a suite of single-turn datasets designed to evaluate diverse reasoning abilities. **MMLU-Pro** (Zheng et al., 2024) extends the original MMLU benchmark with stronger robustness and challenge by filtering data leakage, reducing ambiguity, and introducing more difficult questions across domains such as engineering, philosophy, and economics, making it a more reliable testbed for assessing multi-disciplinary reasoning. **GPQA-Diamond** (Rein et al., 2024) is a graduate-level benchmark featuring expert-written, “Google-proof” questions in physics and related sciences, with its Diamond split being the most challenging and requiring rigorous multi-step reasoning. **AIME-24** and **AIME-25** (HuggingFaceH4, 2024, 2025) consist of Olympiad-style mathematics problems from the 2024 and 2025 American Invitational Mathematics Examinations, testing symbolic manipulation and problem-solving under strict exact-match criteria. Finally, **ToolBench** (Patil et al., 2023) assesses a model’s ability to identify and configure external APIs, reflecting practical tool-use capabilities.

We then evaluate on a suite of multi-turn, goal-oriented benchmarks designed to evaluate memory in embodied and interactive environments. It includes several representative domains: **Alf-World** (Shridhar et al., 2021) for household instruction following, **BabyAI** (Chevalier-Boisvert et al., 2019) for grounded navigation and compositional reasoning, **ScienceWorld** (Wang et al., 2022) for open-ended scientific experimentation, **Jericho** (Hausknecht et al., 2020) for text-based game exploration, and **PDDL** tasks (Yang et al., 2023) for symbolic planning. Together, these environments emphasize long-horizon reasoning, sequential decision-making, and the use of accumulated experience to achieve complex goals.

Together, these datasets form a comprehensive benchmark suite that evaluates factual recall, domain expertise, mathematical reasoning, and procedural memory in interactive settings. This diversity enables a unified evaluation of both static and evolving capabilities, reflecting how LLMs learn, act, and adapt across academic and real-world scenarios.

A.2. Configuration

For efficient retrieval and fair comparison across methods, we utilize the BAAI/bge-base-en-v1.5 (Chen et al., 2023) encoder as the retriever to index both queries and memory items. During inference, the current question is encoded as a query and compared with all stored memory embeddings, retrieving the top- k most relevant items (default $k = 4$) for contextual augmentation. This setting ensures a consistent retrieval budget across all methods. For efficiency, retrieved texts and task inputs are truncated to fit within the same prompt length constraint used by the generation models.

While all baselines adopt the same retrieval configuration, certain methods (e.g., SELF-RAG, REMEM) introduce additional reasoning modules that determine *whether to retrieve* and *what to retrieve* at each step. These adaptive behaviors operate on top of the same retrieval pool to ensure comparability.

Across all experiments, we maintain a unified task sequence ordering within each dataset, ensuring

consistent memory evolution dynamics for all models. Unless otherwise specified, retrieval and generation operate within the same pipeline, and the retrieved items are appended to the prompt following the order of relevance, from most to least similar.

We benchmark a broad range of agents and memory architectures instantiated on two strong **LLM backbones**: the Gemini-2.5 series (Comanici et al., 2025) (FLASH, FLASH-LITE, and PRO) and the Claude family (Anthropic, 2025) (3.5-HAIKU and 3.7-SONNET).

A.3. Evaluation

Evo-Memory evaluates both task performance and memory quality along four key dimensions:

- **Answer accuracy.** Evaluates whether the LLM produces correct outputs across tasks, reflecting its ability to incorporate past experiences into inference.
- **Success rate.** Measures whether the LLM agent successfully completes task goals, indicating its overall effectiveness in interactive or goal-oriented settings.
- **Step efficiency.** Tracks the number of steps required to complete a goal, assessing whether memory usage enables concise and scalable reasoning.
- **Sequence robustness.** Examines whether the LLM maintains consistent knowledge and performance across varying task orders, reflecting its ability to stably reuse prior experiences.

A.4. Methods

We benchmark Evo-Memory with a wide spectrum of agent and memory architectures to study how different designs impact *test-time memory evolution*. All methods are instantiated on two strong **LLM backbones**: Gemini-2.5 (Comanici et al., 2025) and Claude-3.5/3.7 (Anthropic, 2025). Our comparisons isolate the impact of memory architecture and update strategy. Differences in backbone capability are not the focus of the study. We group the evaluated approaches into four major families:

Agent Pipelines without Persistent Memory. ReAct (Yao et al., 2023) serves as a representative reasoning–action pipeline, where memory is limited to the immediate context. It generates interleaved reasoning traces and tool calls but does not explicitly store or evolve information. Amem extends this pipeline with a lightweight agentic memory that caches recent observations and reflections. It provides a minimal form of experience reuse without dedicated search or update policies, forming a bridge between memory-free agents and adaptive memory systems.

Adaptive Agentic Memory Methods. This group focuses on adaptive retrieval and self-evolving memory. SelfRAG (Asai et al., 2023) integrates dynamic retrieval and reflection to adaptively ground reasoning in prior contexts. MemOS (Li et al., 2025), Mem0 (Chhikara et al., 2025), and LangMem (LangChain contributors, 2025) implement structured, agent-level memory systems that support read, write, and update operations. Within our unified interface, retrieval corresponds to the *search* stage and updates correspond to *evolve*. These methods represent adaptive long-term agents capable of continual refinement.

Memory-Based Agents for Procedural Memory. Dynamic Cheatsheet (DC) (Suzgun et al., 2025) and Agent Workflow Memory (AWM) (Wang et al., 2024) emphasize the reuse of procedural knowledge, encoding “how-to” information rather than static facts. We evaluate two DC variants, DC-RS (retrieval-based) and DC-Cu (curated), to analyze how workflow induction and update mechanisms

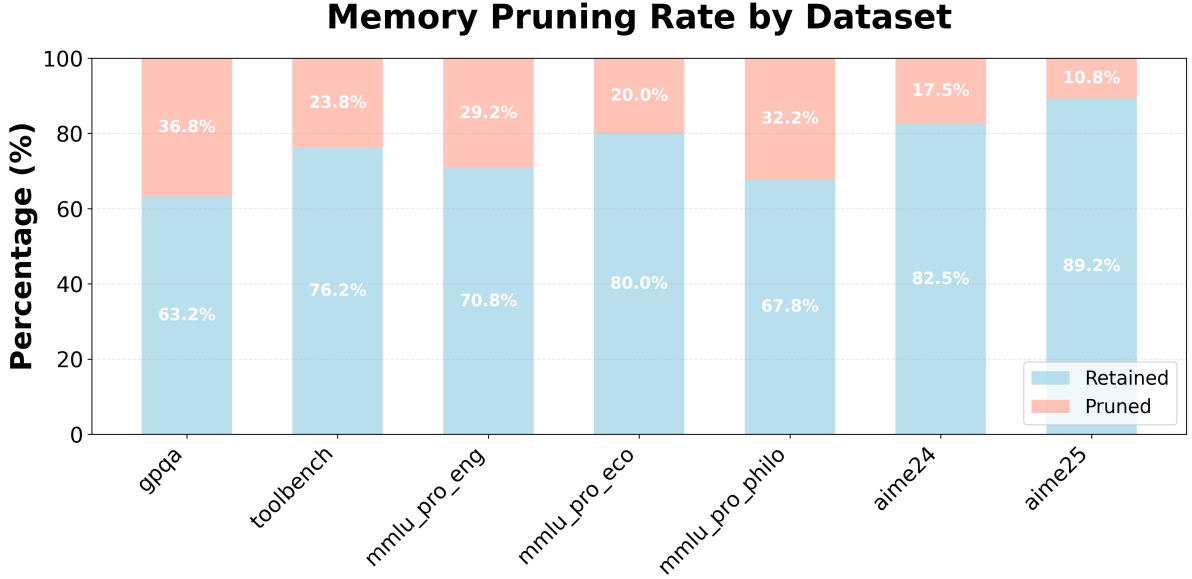


Figure 7 | Memory pruning rates by dataset. Retained (blue) and pruned (coral) memory proportions show varying selectivity across benchmarks.

influence stability and transfer. These methods test the potential of procedural memory as reusable strategy repositories.

Proposed: Evolving Memory Framework. **ExpRecent** maintains condensed episodic traces of recent task trajectories, while our **ExpRAG** family integrates the principles of retrieval-augmented reasoning with explicit *test-time evolution*. **ReMem** applies iterative reflection and synthesis to refine memory embeddings over time. Together, these methods instantiate Evo-Memory’s design philosophy, treating reasoning, acting, and memory refinement as interleaved processes that co-adapt during deployment, enabling continual self-improvement and more human-like adaptation.

B. Experiments

We provide more experiments in the following.

B.1. Additional Experiments

We further validate our findings through extensive benchmarking across multiple model families (Gemini-2.5-Flash-Lite, Claude-3.5-Haiku) and diverse datasets, as shown in Tables 5 and 6. The performance trends remain consistent across all settings. On both multi-turn embodied reasoning tasks (Alf World, BabyAI, PDDL, ScienceWorld) and single-turn reasoning tasks (AIME-24/25, GPQA, MMLU-Pro, ToolBench), ReMem consistently outperforms conventional baselines and adaptive retrieval methods across model backbones. These results confirm that the advantages of evolving-memory architectures are model-agnostic, highlighting continual task-level reflection as a general mechanism for improving adaptability in problem-solving.

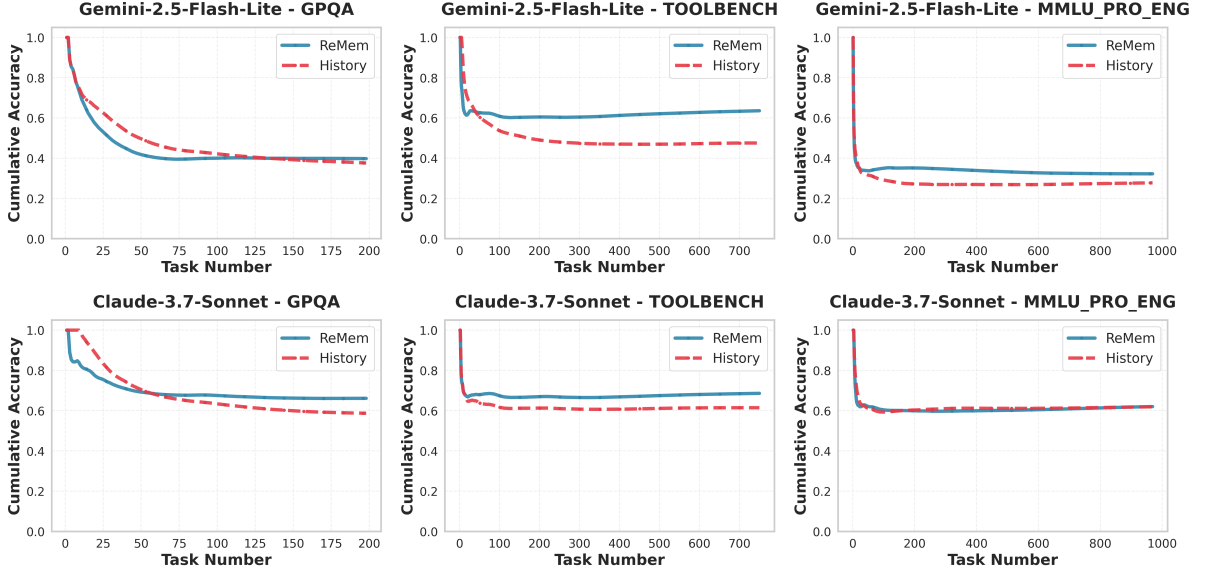


Figure 8 | Cumulative accuracy comparison across model variants and benchmarks. ReMem (solid blue) demonstrates consistent improvements over the History baseline (dashed red) across Gemini-2.5-Flash-Lite and Claude-3.7-Sonnet models on GPQA, TOOLBENCH, and MMLU_PRO_ENG datasets. The curves show learning trends as task instances accumulate, with ReMem achieving faster convergence and higher final accuracy.

B.2. Additional Analysis of Memory Pruning

Figure 7 shows memory pruning rates across datasets, revealing varying selectivity in memory retention. The pruning ratios differ substantially across benchmarks, which appears related to task diversity and domain coverage. Datasets with broader domain coverage such as GPQA, which encompasses diverse problem types across engineering, physics, and other domains, exhibit higher pruning rates (36.8%), suggesting that more memories are deemed redundant across heterogeneous tasks. In contrast, datasets with more concentrated problem types like AIME show lower pruning rates (17.5% and 10.8% respectively), indicating that memories remain more relevant due to higher task similarity. This pattern suggests that the pruning mechanism effectively identifies and discards domain-irrelevant experiences, though the precise relationship between task diversity and memory selectivity warrants further investigation.

B.3. Additional Comparative Curves on Single-turn Tasks

Figure 8 presents cumulative accuracy curves comparing ReMem with the baseline across single-turn reasoning benchmarks and model variants. As task instances accumulate, ReMem shows consistent improvement on GPQA, ToolBench, and MMLU-PRO (Engineer) for both Gemini-2.5-Flash-Lite and Claude-3.7-Sonnet. Similar to the multi-turn results, History performs comparably at the beginning due to the cold-start phase, but ReMem quickly surpasses it as more tasks are processed, indicating the cumulative advantage of continual task-level adaptation.

C. Prompts

Memory Prompt Template for Multi-turn Dataset

=====

ENVIRONMENT INSTRUCTIONS

=====

[Detailed task environment description and rules]
Example: Go to kitchen, pick up apple, put it in bag, etc.

=====

EXAMPLE DEMONSTRATIONS

=====

[Static few-shot examples]
Example 1: Goal: ... | Action: ... | Observation: ...
Example 2: Goal: ... | Action: ... | Observation: ...

=====

RELEVANT EXPERIENCE FROM SIMILAR TASKS

=====

[Experience #1]
Goal: [similar goal]
Trajectory: [action sequence]
Correctness: [success/failure]
[Experience #2, #3, ...]

=====

YOUR CURRENT TASK

=====

Goal: *[specific task goal]*
Help: type 'check valid actions' if action fails
Help: type 'inventory' to check items

=====

RECENT HISTORY

=====

Observation: *[initial environment state]*
 Action: *[previous action]*
 Observation: *[result of previous action]*
 Action: *[previous action]*
 Observation: *[current state]*

=====

OUTPUT FORMAT

=====

You MUST respond in EXACTLY ONE of these formats:

Format 1 - Prune experiences:

Think-Prune: *<IDs>*
 Remove unhelpful experiences from 'RELEVANT EXPERIENCE' section (e.g., "1,3" or "2-4")

Format 2 - Internal reasoning:

Think: *<your reasoning>*
 Free-form explanation of your next step

Format 3 - Execute action:

Action: *<exact command>*
 Must be valid command from ENVIRONMENT INSTRUCTIONS with exact names from RECENT HISTORY

Memory Prompt Template for Single-turn Dataset

You are a helpful assistant with access to LOCAL EXPERIENCE MEMORY. Each memory may contain past experience, rationales, domains, and skills. Below are some retrieved LOCAL EXPERIENCE MEMORIES:

[Retrieved/synthesized memories]

Now solve the following problem.

Question: *[Your question here]*

Provide your output in the following format:

- **Rationale:** your short reasoning, may cite memory if useful
- **Final Answer:** your final answer

D. Limitations

While Evo-Memory offers a comprehensive evaluation of self-evolving memory, several practical constraints remain. Due to budget and API limits, we focus on a selected set of strong LLMs rather than exhaustively covering all available models. Additional evaluations on open-weight or multilingual models could further validate the generality of our findings. Moreover, our benchmark primarily emphasizes textual and goal-oriented tasks; extending it to richer multimodal or real-world environments would provide a more complete picture of continual memory evolution. Despite these limitations, the current study already spans diverse domains, tasks, and architectures, offering a solid foundation for future extensions.

E. Use of Large Language Models

During the preparation of this paper, we made limited and controlled use of large language models (LLMs), specifically ChatGPT, as an auxiliary writing aid. The LLM was used only for stylistic refinement, including improvements in clarity, grammar, and readability of text originally drafted by the authors. All scientific ideas, analyses, experiments, and conclusions were fully developed, written, and verified by the authors. Thus, LLMs were employed solely as a language-editing tool, without contributing to the intellectual or scientific content of the work.

LLM Backbone	Method	AlfWorld		BabyAI		PDDL		ScienceWorld		Avg.	
		S	P	S	P	S	P	S	P	S	P
Gemini 2.5 Flash	Baseline	0.12	0.34	0.61	0.71	0.12	0.20	0.24	0.59	0.27	0.46
	History	0.28	0.60	0.52	0.64	0.08	0.15	0.31	0.71	0.30	0.53
	ReAct	0.24	0.56	0.48	0.63	0.22	0.33	0.34	0.71	0.32	0.56
	Amem	0.25	0.59	0.53	0.64	0.10	0.16	0.36	0.74	0.31	0.53
	SelfRAG	0.25	0.59	0.52	0.65	0.08	0.16	0.34	0.74	0.30	0.54
	Mem0	0.27	0.61	0.54	0.66	0.10	0.19	0.32	0.70	0.31	0.54
	DC-Cu	0.25	0.59	0.53	0.64	0.08	0.17	0.29	0.71	0.29	0.53
	DC-RS	0.27	0.60	0.53	0.66	0.07	0.15	0.33	0.73	0.30	0.54
	AWM	0.26	0.59	0.52	0.64	0.08	0.16	0.33	0.73	0.30	0.53
	ExpRecent	0.37	0.65	0.53	0.64	0.13	0.22	0.53	0.83	0.39	0.59
	ExpRAG	0.59	0.79	0.56	0.65	0.17	0.27	0.53	0.81	0.46	0.63
	ReMem	0.66	0.81	0.53	0.61	0.22	0.33	0.58	0.81	0.50	0.64
Gemini 2.5 Pro	Baseline	0.04	0.39	0.37	0.47	0.20	0.33	0.22	0.60	0.21	0.45
	History	0.19	0.52	0.40	0.48	0.25	0.38	0.60	0.84	0.36	0.56
	ReAct	0.02	0.26	0.43	0.55	0.13	0.22	0.30	0.68	0.22	0.43
	Amem	0.16	0.50	0.42	0.49	0.23	0.38	0.59	0.85	0.35	0.56
	SelfRAG	0.16	0.49	0.43	0.50	0.22	0.35	0.57	0.84	0.34	0.55
	Mem0	0.16	0.49	0.41	0.49	0.22	0.37	0.53	0.81	0.33	0.54
	DC-Cu	0.17	0.50	0.40	0.47	0.28	0.40	0.53	0.82	0.35	0.55
	DC-RS	0.18	0.51	0.42	0.50	0.27	0.40	0.59	0.85	0.37	0.57
	AWM	0.20	0.52	0.38	0.46	0.20	0.37	0.57	0.83	0.34	0.54
	ExpRecent	0.36	0.61	0.54	0.64	0.35	0.47	0.69	0.89	0.49	0.64
	ExpRAG	0.38	0.64	0.46	0.53	0.28	0.43	0.61	0.84	0.43	0.61
	ReMem	0.51	0.70	0.56	0.64	0.25	0.38	0.66	0.86	0.50	0.65
Claude 3.7 Sonnet	Baseline	0.18	0.49	0.51	0.66	0.17	0.39	0.10	0.53	0.24	0.52
	History	0.50	0.73	0.48	0.66	0.65	0.85	0.32	0.74	0.49	0.74
	ReAct	0.51	0.75	0.57	0.72	0.75	0.91	0.44	0.77	0.57	0.79
	Amem	0.48	0.73	0.46	0.64	0.62	0.84	0.33	0.73	0.47	0.73
	SelfRAG	0.52	0.75	0.46	0.64	0.65	0.84	0.31	0.74	0.49	0.74
	Mem0	0.51	0.74	0.48	0.66	0.65	0.84	0.37	0.76	0.50	0.75
	DC-Cu	0.50	0.74	0.50	0.67	0.62	0.84	0.33	0.75	0.49	0.75
	DC-RS	0.50	0.74	0.52	0.68	0.62	0.84	0.34	0.74	0.50	0.75
	AWM	0.49	0.73	0.53	0.68	0.60	0.82	0.34	0.74	0.49	0.74
	ExpRecent	0.66	0.83	0.63	0.73	0.53	0.76	0.49	0.82	0.58	0.79
	ExpRAG	0.74	0.89	0.62	0.72	0.72	0.89	0.46	0.76	0.63	0.82
	ReMem	0.92	0.96	0.73	0.83	0.83	0.95	0.62	0.89	0.78	0.91
Claude 3.5 Haiku	Baseline	0.11	0.33	0.38	0.52	0.15	0.32	0.08	0.37	0.18	0.39
	History	0.28	0.58	0.38	0.57	0.18	0.38	0.12	0.49	0.24	0.51
	ReAct	0.24	0.58	0.35	0.52	0.32	0.53	0.16	0.55	0.27	0.55
	Amem	0.24	0.55	0.37	0.58	0.17	0.35	0.12	0.45	0.23	0.48
	SelfRAG	0.26	0.58	0.38	0.59	0.22	0.37	0.14	0.49	0.25	0.51
	Mem0	0.27	0.56	0.37	0.57	0.17	0.37	0.08	0.45	0.22	0.49
	DC-Cu	0.24	0.55	0.37	0.58	0.17	0.37	0.12	0.45	0.23	0.49
	DC-RS	0.24	0.55	0.37	0.58	0.17	0.37	0.12	0.45	0.23	0.49
	AWM	0.24	0.55	0.37	0.58	0.17	0.37	0.12	0.45	0.23	0.49
	ExpRecent	0.48	0.65	0.40	0.57	0.15	0.32	0.32	0.64	0.34	0.55
	ExpRAG	0.65	0.74	0.54	0.64	0.43	0.61	0.42	0.68	0.51	0.67
	ReMem	0.69	0.80	0.49	0.60	0.43	0.61	0.44	0.75	0.51	0.69

Table 5 | Cross-environment results across four embodied reasoning benchmarks (Alf World, BabyAI, PDDL, ScienceWorld). Each dataset reports success (S) and progress (P) rates. Bold indicates the best (including ties) per column. The last two columns show averaged S and P across datasets.

LLM Backbone	Method	Exact Match \uparrow						API / Acc. \uparrow	
		AIME24	AIME25	GPQA	MMLU-Pro (Eco.)	MMLU-Pro (Eng.)	MMLU-Pro (Philo.)	ToolBench	Avg. \uparrow
Claude 3.5	Baseline	—	—	0.36	0.68	0.42	0.55	0.81/0.64	0.38
	History	—	—	0.37	0.70	0.43	0.55	0.81/0.63	0.38
	ReAct	—	—	0.35	0.69	0.43	0.54	0.81/0.64	0.38
	Amem	—	—	0.34	0.69	0.43	0.53	0.82/0.63	0.37
	SelfRAG	—	—	0.36	0.70	0.44	0.56	0.83/0.65	0.39
	MemOS	—	—	0.37	0.70	0.42	0.55	0.81/0.64	0.38
	Mem0	—	—	0.36	0.70	0.42	0.55	0.82/0.64	0.38
	LangMem	—	—	0.51	0.78	0.46	0.61	0.81/0.63	0.43
	DC-RS	—	—	0.36	0.68	0.41	0.56	0.82/0.64	0.38
	AWM	—	—	0.33	0.67	0.42	0.53	0.81/0.63	0.37
	DC-Cu	—	—	0.33	0.65	0.39	0.54	0.82/0.63	0.36
	ExpRecent	—	—	0.42	0.69	0.46	0.59	0.85/0.63	0.40
	ExpRAG	—	—	0.40	0.73	0.49	0.61	0.87/0.67	0.41
	ReMem	—	—	0.39	0.71	0.47	0.62	0.87/0.68	0.41
Claude 3.7 Sonnet	Baseline	0.17	0.13	0.55	0.84	0.63	0.78	0.76/0.62	0.54
	History	0.13	0.23	0.56	0.85	0.64	0.78	0.76/0.61	0.55
	ReAct	0.17	0.10	0.57	0.84	0.63	0.76	0.76/0.61	0.54
	Amem	0.27	0.17	0.54	0.83	0.63	0.79	0.77/0.63	0.56
	SelfRAG	0.20	0.10	0.58	0.84	0.65	0.77	0.77/0.63	0.55
	MemOS	0.17	0.20	0.55	0.84	0.64	0.76	0.76/0.62	0.55
	Mem0	0.20	0.13	0.58	0.84	0.62	0.77	0.76/0.61	0.55
	LangMem	0.10	0.13	0.53	0.77	0.56	0.66	0.77/0.63	0.49
	DC-RS	0.20	0.20	0.62	0.79	0.52	0.60	0.77/0.62	0.52
	AWM	0.03	0.03	0.53	0.80	0.56	0.72	0.76/0.62	0.48
	DC-Cu	0.17	0.23	0.57	0.79	0.52	0.65	0.77/0.62	0.52
	ExpRecent	0.13	0.20	0.61	0.86	0.63	0.78	0.82/0.66	0.56
	ExpRAG	0.17	0.17	0.70	0.85	0.67	0.80	0.88/0.72	0.59
	ReMem	0.13	0.13	0.67	0.86	0.65	0.80	0.87/0.71	0.58
Gemini 2.5 Flash	Baseline	0.47	0.47	0.48	0.83	0.46	0.75	0.71/0.61	0.59
	History	0.60	0.47	0.43	0.84	0.42	0.78	0.31/0.26	0.55
	ReAct	0.30	0.27	0.05	0.64	0.16	0.54	0.64/0.57	0.37
	Amem	0.70	0.57	0.52	0.83	0.42	0.72	0.72/0.60	0.63
	SelfRAG	0.50	0.47	0.46	0.83	0.45	0.75	0.72/0.61	0.59
	MemOS	0.47	0.47	0.50	0.82	0.46	0.75	0.71/0.61	0.59
	Mem0	0.50	0.47	0.45	0.83	0.46	0.74	0.71/0.61	0.59
	LangMem	0.43	0.50	0.53	0.79	0.39	0.71	0.68/0.57	0.57
	DC-RS	0.53	0.37	0.48	0.80	0.42	0.69	0.68/0.57	0.56
	DC-Cu	0.60	0.40	0.48	0.79	0.44	0.69	0.70/0.59	0.58
	AWM	0.50	0.37	0.49	0.79	0.43	0.72	0.71/0.59	0.56
	ExpRecent	0.47	0.47	0.42	0.83	0.39	0.75	0.78/0.66	0.58
	ExpRAG	0.43	0.47	0.42	0.83	0.43	0.78	0.87/0.73	0.60
	ReMem	0.60	0.53	0.51	0.85	0.46	0.79	0.85/0.71	0.65
Gemini 2.5 Flash-Lite	Baseline	0.53	0.43	0.37	0.73	0.34	0.60	0.78/0.61	0.58
	History	0.40	0.33	0.31	0.74	0.30	0.59	0.58/0.47	0.49
	ReAct	0.53	0.33	0.34	0.61	0.20	0.48	0.73/0.56	0.50
	Amem	0.40	0.33	0.33	0.72	0.34	0.59	0.77/0.61	0.54
	SelfRAG	0.57	0.37	0.37	0.73	0.32	0.62	0.81/0.63	0.57
	MemOS	0.53	0.43	0.37	0.73	0.34	0.60	0.78/0.61	0.58
	Mem0	0.53	0.43	0.37	0.73	0.34	0.60	0.78/0.61	0.58
	LangMem	0.40	0.37	0.48	0.59	0.24	0.56	0.33/0.26	0.43
	DC-RS	0.53	0.43	0.34	0.59	0.18	0.36	0.73/0.56	0.49
	AWM	0.03	0.03	0.35	0.61	0.20	0.48	0.77/0.60	0.44
	DC-Cu	0.53	0.43	0.33	0.55	0.16	0.32	0.71/0.56	0.47
	ExpRecent	0.57	0.33	0.35	0.76	0.29	0.62	0.82/0.65	0.58
	ExpRAG	0.47	0.37	0.38	0.79	0.32	0.66	0.87/0.68	0.61
	ReMem	0.57	0.33	0.38	0.77	0.34	0.65	0.86/0.67	0.61

Table 6 | Cross-dataset results of diverse memory architectures across models. Categories are separated by horizontal rules; results (Exact Match \uparrow and API/Acc \uparrow) compare zero-shot, agentic, adaptive, procedural, and proposed memory methods. Dashes (—) indicate methods with poor or unreliable performance, which are therefore omitted.