

Mem^p: Exploring Agent Procedural Memory

Anonymous ACL submission

Abstract

Large Language Models (LLMs) based agents excel at diverse tasks, yet they suffer from brittle procedural memory that is manually engineered or entangled in static parameters. In this work, we investigate strategies to endow agents with a *learnable, updatable, and life-long* procedural memory. We propose Mem^p that distills past agent trajectories into both fine-grained, step-by-step instructions and higher-level, script-like abstractions, and explore the impact of different strategies for *Build, Retrieval, and Update* of procedural memory. Coupled with a dynamic regimen that continuously updates, corrects, and deprecates its contents, this repository evolves in lockstep with new experience. Empirical evaluation on TravelPlanner and ALFWorld shows that as the memory repository is refined, agents achieve steadily higher success rates and greater efficiency on analogous tasks. Moreover, procedural memory built from a stronger model retains its value: migrating the procedural memory to a weaker model can also yield substantial performance gains.

1 Introduction

As large language models (LLMs) grow ever more powerful, LLM-based agents augmented by their own reasoning and external tools are taking on increasingly sophisticated works (Zhao et al., 2023; Wang et al., 2024a; Xi et al., 2025; Qiao et al., 2023). No longer mere assistants, these agents now trawl the web for elusive insights and weave them into comprehensive, publication ready reports, like Deep Research (OpenAI, 2025; x.ai, 2025) and WebDancer (Wu et al., 2025a). Moreover, they can handle complex data analyses (Lan et al., 2025; Ifargan et al., 2025; Ou et al., 2025), navigate multi-step GUI workflows (Luo et al., 2025; Qin et al., 2025), and sustain long-horizon, tool-rich interactions (Yao et al., 2025; Barres et al., 2025; Chen et al., 2025; Fang et al., 2025; Gur et al., 2023)

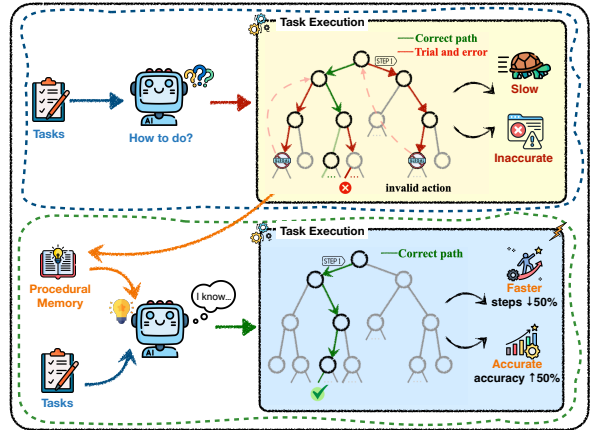


Figure 1: With **procedural memory**, agents can improve both the success rate (accuracy \uparrow) and execution efficiency (steps \downarrow) when solving similar tasks.

with precision. Yet executing such intricate, long-horizon tasks demands dozens of steps and protracted runtimes. Along the way, unpredictable external events—network glitches, UI changes, shifting data schemas—can derail the entire process. Restarting from scratch every time is a punishing ordeal for present-day agents. Beneath their surface diversity, many complex tasks share deep structural commonalities and a similar environment. Instead of starting fresh each time, an agent should extract its experience from past successes. By turning earlier trajectories into reusable templates like patterns of reasoning, tool sequences, and recovery tactics, it can progress step by step, learning from every failure and success, until even the most convoluted missions become routine.

The capacity to distill, chronicle, and re-apply lessons from one’s own experiential trajectory is the bedrock of human learning and the pivotal gateway through which an agent ascends toward self-directed refinement (Liu et al., 2025a; Summers et al., 2023a; Li et al., 2023). Procedural memory (Gupta and Cohen, 2002; Cohen and Squire, 1980) silently compiles habitual skills into exe-

067 cutable subroutines, enabling unconscious, fluent
068 action. While contemporary agents built on LLMs
069 can compose short action plans or call external
070 tools, their procedural knowledge is either hand-
071 crafted, stored as brittle prompt templates, or im-
072 plicitly entangled in model parameters that are ex-
073 pensive to update. Existing memory-augmented
074 frameworks such as LangGraph (Mavroudis, 2024),
075 AutoGPT (Yang et al., 2023), or agent cognitive ar-
076 chitectures like Memory Bank (Zhong et al., 2024a;
077 Sumers et al., 2023b) and Soar (Laird, 2022) pro-
078 vide coarse abstractions (buffers, rule chunks, pro-
079 duction systems) but leave the optimization of pro-
080 cedural memory life-cycle operations about how
081 skills are built, indexed, patched, and eventually
082 pruned, largely unexamined. Consequently, there
083 is no reliable way to measure how efficiently an
084 agent grows its procedural skills or to ensure new
085 experiences improve rather than harm performance.

086 To close this gap, we present *Mem^P*, a task-
087 agnostic framework that treats procedural memory
088 as a first-class optimization object. The core explo-
089 ration of *Mem^P* lies in how different strategies for
090 memory construction, retrieval, and updating af-
091 fect overall performance. During the construction
092 phase, we follow the majority of traditional mem-
093 ory architectures and agent-based memory designs
094 by leveraging either the full historical trajectory
095 or explicit guidelines to guide the process. In the
096 retrieval phase, we experiment with various key-
097 building strategies—such as query-vector match-
098 ing and keyword-vector matching—to investigate
099 how procedural memory can be constructed more
100 precisely. Unlike prior memory mechanisms or
101 learning from experience, *Mem^P* introduces di-
102 verse procedural-memory update strategies: In the
103 realm of agents, memory updating is crucial for
104 agents to adapt to dynamic environments. By incor-
105 porating diverse strategies like ordinary addition,
106 validation filtering, reflection, and dynamic discard-
107 ing, agents can efficiently manage their knowledge
108 base. This ensures they stay updated with new
109 information, discard outdated data, and optimize
110 memory resources. Such strategies enhance learn-
111 ing efficiency, improve decision-making quality,
112 and boost adaptability, allowing agents to perform
113 optimally in various tasks and scenarios.

114 We instantiate *Mem^P* on top of strong LLMs
115 (GPT-4o and Claude, Qwen) and evaluate on two
116 diverse domains: long-horizon housework ALF-
117 World (Shridhar et al., 2021) and long-term in-
118 formation seeking task TravelPlanner (Xie et al.,

2024). On two benchmark datasets that rigorously
evaluate agent capabilities, we demonstrate that
constructing and retrieving procedural memory dur-
ing training empowers an agent to distill and reuse
its prior experience. When this memory is ex-
ploited at test time, the agent’s task accuracy rises,
and compared with tackling each instance in iso-
lation, it eliminates most fruitless exploration on
unfamiliar tasks, yielding substantial reductions in
both step count and token consumption. Further, by
equipping the agent with a set of memory-update
mechanisms, we allow it to build and refine its
procedural memory while acting in the test envi-
ronment. This endows the agent with a continual,
almost linear mastery of the task. Extensive abla-
tion studies reveal that procedural memory not only
scales gracefully with increasing task complexity
but also transfers effectively to new, related tasks,
demonstrating its adaptability across scenarios.

2 Related Works

Memory in Language Agents. Memory is a
foundational component in language agents, en-
abling them to retain and utilize past information
across multiple timescales, including short-term,
episodic, and long-term memory, to enhance their
performance and adaptability (Zhou et al., 2023,
2024; Zhang et al., 2024; Liu et al., 2025a; Li
et al., 2025). These systems aim to mimic as-
pects of human memory to improve coherence, per-
sonalization, and learning capabilities (Chhikara
et al., 2025; Wu et al., 2025b; Xia et al., 2025).
Current approaches include end-to-end memory
systems (Yu et al., 2025; Zhou et al., 2025), ex-
ternal memory systems (Chhikara et al., 2025;
Zhong et al., 2024b), and hierarchical memory
structures (Hu et al., 2024a; Xu et al., 2025). These
methods involve encoding and storing information
in various formats, using retrieval mechanisms like
vector embeddings and semantic search, and imple-
menting memory updating and forgetting strategies
to maintain relevance and efficiency. Despite its im-
portance, memory in multi-turn agent interactions
remains underexplored, and enabling agents to ef-
fectively learn and utilize memory across trajec-
tories poses a significant challenge. Procedural mem-
ory is a type of long-term memory that involves the
retention of procedures and skills, such as typing
or riding a bike, which are performed automati-
cally without conscious thought. The agent utilizes
procedural memory to internalize and automate

repetitive tasks, decision-making processes, and interaction patterns, leading to more efficient and context-aware responses over time. Although there have been several works, such as Voyager (Wang et al., 2023), AWM (Wang et al., 2024b), and AutoManual (Chen et al., 2024), that utilize procedural memory to enhance agents’ capabilities on similar tasks, there still lacks a systematic analysis on how to construct, retrieve, and update such procedural memory like (Wu et al., 2024). Therefore, our work mainly focuses on exploring how to build an effective procedural memory system for agents performing cross-trajectory tasks.

Learning from Experience. LLM-based Agent learning from experience involves intelligence continuously improving their decision-making capabilities through interaction with environments and utilization of past experiences (Tan et al., 2025; Tang et al., 2025; Zhou et al., 2025; Qiao et al., 2025; Su et al., 2025; Wang et al., 2024b). This approach is crucial for developing adaptive and intelligent agents capable of handling dynamic real-world scenarios, as it allows them to optimize behaviors, reduce manual programming needs, and enhance performance across various tasks (Zheng et al.; Liu et al., 2025c; Wang et al., 2025). Agents typically employ mechanisms such as reinforcement learning (Lu et al., 2025; Dong et al., 2025), experience replay (Feng et al., 2025; Liu et al., 2025b), imitation learning (Sun et al., 2024; Yang et al., 2024b), memory management (Hou et al., 2024; Hu et al., 2024b), and multi-agent learning to achieve this. However, current methods face limitations including low sample efficiency, poor generalization across tasks, catastrophic forgetting when learning new information, and there are very few features for memory update. The key distinction of our work lies in systematically investigating optimal strategies for construction, retrieval, and update modules of an agent’s procedural knowledge. During the update phase, we enhance the agent’s capabilities by maintaining an editable repository of procedural knowledge. Additionally, collecting high-quality training data can be challenging and may introduce biases.

3 Preliminary

When an agent influences its external environment by invoking external tools or executing prescribed actions, and iteratively refines its behavior over multiple rounds to accomplish a complex multi-

step objective, this paradigm can be modeled as a Markov Decision Process (MDP). (Puterman, 1990) Under this view, at each discrete time step t , the agent, situated in state $s_t \in S$, chooses an action $a_t \in A$, according to its policy $\pi(a_t|s_t)$, where A is the action space of the task. The environment then transitions to a new state $s_{t+1} \in S$ and emits an observation O_t . Consequently, the entire interaction trajectory may be compactly expressed as:

$$\tau = (s_0, a_0, o_1, s_1, a_1, o_2, \dots, s_T), \quad (1)$$

where τ is the complete exploration trajectory of this task. Moreover, a reward function R will evaluate the task’s completion r within this environment env by assigning a score based on the final state s_T or the entire trajectory τ .

$$r = R(env, s_T, \tau) \in [0, 1] \quad (2)$$

Although approaches resembling Markov Decision Processes inevitably contain erroneous actions and exploratory attempts, the contextual information they generate becomes valuable for decision-making as the model’s reasoning and reflective capabilities improve. Nevertheless, this benefit comes at a high test-time cost—both in time and in token consumption. When facing an entirely new and complex environment, many actions (or tokens) are spent simply understanding the environment and the task itself. This leads to redundancy when similar tasks are executed within the same environment: the agent has already acquired partial procedural knowledge about the environment or task during earlier episodes, yet fails to transfer that knowledge effectively to subsequent tasks.

By shifting from parallel task completion to sequential task completion, the agent can learn and distill experience from earlier tasks, thereby reducing repetitive exploration. Inspired by human procedural memory, we propose to equip the agent with a procedural memory module. This module transforms the conventional policy $\pi(a_t|s_t)$ into $\pi_{m^p}(a_t|s_t)$, where m^p is the agent’s learned procedural memory.

3.1 Agent Procedural Memory

Procedural memory is the type of long-term memory responsible for knowing how to perform tasks and skills, such as typing or riding a bike. By mastering this type of procedural memory, humans avoid the need to relearn the process each time. For

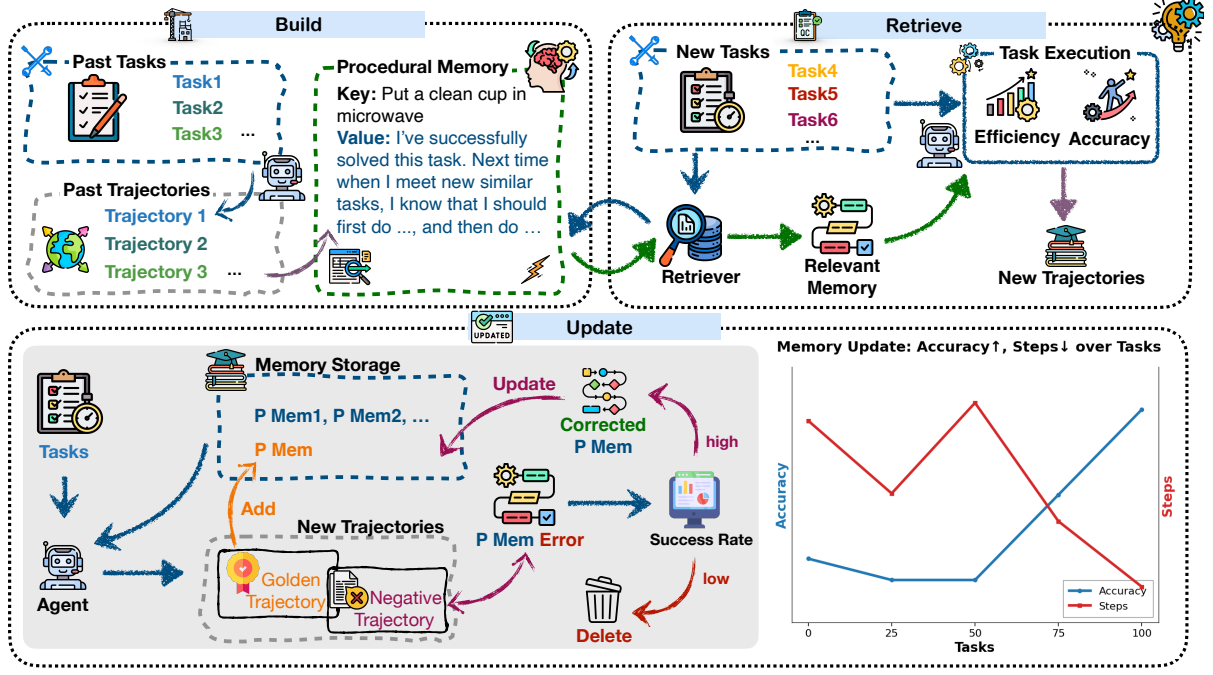


Figure 2: The procedural memory framework consists of **Build**, **Retrieve**, and **Update**, which respectively involve encoding stored procedural memory, forming new procedural memories, and modifying existing ones in light of new experiences.

an agent, that is, for a task trajectory τ and its reward r , a memory m^p is constructed by a builder B , thereby achieving the acquisition of memory, namely

$$Mem = \sum_{t=1}^T m^{pt}, \text{ where } m^{pt} = B(\tau_t, r_t) \quad (3)$$

where Mem is the procedural memory library acquired by the agent over the T tasks. After constructing the procedural memory library, when facing a new task t_{new} , we need a good procedural memory retriever to recall a memory that fits t_{new} . Generally speaking, we would choose the task $t \in T$ that is most similar to t_{new} , because similar experiences are more helpful for the agent to complete the new task.

$$m_{retrieved} = \arg \max_{m^{pi} \in Mem} S(t_{new}, t_i) \quad (4)$$

As we use cosine similarity for the vector embedding model ϕ of the task in the experiment, the retrieval process becomes:

$$m_{retrieved} = \arg \max_{m^{pi} \in Mem} \frac{\phi(t_{new}) \cdot \phi(t_i)}{\|\phi(t_{new})\| \|\phi(t_i)\|}. \quad (5)$$

Moreover, as the number of completed tasks continuously increases, simply augmenting the agent's

procedural memory is inconsistent with common sense. A well-designed procedural memory system should have a reasonable update mechanism—that is, it should dynamically perform addition, deletion, modification, and retrieval based on the task execution context.

Let $M(t)$ denote the agent's procedural memory at time t , and τ_t represent the set of tasks completed up to time t . Then, the update mechanism can be modeled as a function U that takes the current procedural memory and task execution feedback to produce the updated memory:

$$M(t+1) = U(M(t), E(t), \tau_t), \quad (6)$$

where $E(t)$ encapsulates the execution feedback (e.g., success, failure, performance metrics). A more sophisticated implementation of U could be represented as:

$$U = Add(M_{new}) \ominus Del(M_{obs}) \oplus Update(M_{est}) \quad (7)$$

where M_{new} represents new procedural memory to be added; M_{obs} indicates procedural memory to be removed, M_{est} are tasks to be updated based on execution feedback $E(t)$. This comprehensive formula captures the essential add, delete, and modify operations within the update mechanism.

4 Experiment

In this section, we will introduce the Procedural Memory framework in detail (Figure 2), covering the storage, retrieval, and update modules of memory, as well as analyzing which strategies perform better within each module.

4.1 Experimental Settings

Datasets. For our experiments, we adapt TravelPlanner (Xie et al., 2024) and ALFWorld (Shridhar et al., 2021) benchmarks. TravelPlanner is a benchmark designed to evaluate agents’ ability to use tools and perform complex planning under intricate constraints. In contrast, ALFWorld comprises household tasks. In each interaction round, the agent outputs an action, and the environment responds with textual feedback describing the resulting state. This process repeats for multiple turns until the task is completed or the maximum number of rounds is reached. ALFWorld includes test split to evaluate the agent’s generalization ability. Detailed evaluation methods are shown in Appendix C.

Backbones. In our experiments, we benchmarked our procedural memory on three base models. Specifically, we adopt the two proprietary frontier models that have consistently dominated public leaderboards: OpenAI’s GPT-4o (OpenAI, 2022) and Anthropic’s Claude (Anthropic, 2022), and complement them with the open-sourced Qwen2.5-72B-Instruct (Yang et al., 2024a). The first two provide state-of-the-art closed-source performance, while the third allows us to verify that our findings generalize beyond proprietary systems and remain valid in the open-source regime.

4.2 Memory Storage & Retrieval

Procedural knowledge is typically stored in two main formats: (1) trajectories are kept verbatim, round by round, in memory, or (2) high-level abstractions are extracted from these trajectories and then stored. Once a similar procedural memory is retrieved, it is appended to the task as part of the context, serving as prior knowledge to assist the model in completing the task.

Inspired by this, we designed the following experimental conditions:

- **No Memory:** The model tackles the assigned task in a ReAct fashion without any external memory.

- **Trajectory:** We first filter the gold trajectories from the training set and store them. At inference time, the system retrieves the top-k trajectories whose query vectors are most similar to the current task’s vector, supplying them as procedural memories before execution.

- **Script:** The model analyzes and summarizes the gold trajectories from the training set, distilling them into abstract procedural knowledge that is provided as a prompt before each task.

- **Proceduralization:** This condition combines the full retrieved trajectories with the high-level script generated by the model, integrating both concrete examples and abstract guidance as the procedural memory.

As shown in Table 1, all memory construction methods outperform the no-memory baseline, achieving higher scores on both datasets while also reducing the number of steps required. This indicates that procedural memory built during training is beneficial for directly applying tasks during testing. Furthermore, we observe that the approach of abstracting trajectories into scripts during training yields relatively better performance on the ALFWorld test set compared to the dev set. Conversely, trajectories that utilize complete execution traces as procedural memory achieve higher scores on the dev set, suggesting that scripts are more capable of generalizing to different test tasks, while trajectories are better suited for scenarios involving tasks similar to those already completed. By combining procedure knowledge from both methods of employing abstracted guidelines along with concrete execution trajectories, we attain the optimal performance.

After converting a set of completed trajectories into procedural memory, the next critical challenge is to retrieve the most accurate and relevant procedural knowledge when a new task arrives. We have designed several different key construction methods for memory storage to facilitate subsequent vector-based matching and retrieval:

- **Random Sample:** Does not utilize keys for vector retrieval; instead, randomly extracts a few memories from procedural memory.
- **Query:** Employ query description as the key for storage, leveraging the semantic similarity

Model	Granularity	TravelPlanner			ALFWorld		
		#CS \uparrow	#HC \uparrow	Steps \downarrow	Dev \uparrow	Test \uparrow	Steps \downarrow
GPT-4o	No Memory	71.93	12.88	17.84	39.28	42.14	23.76
	Script	72.08	5.50	15.79	66.67	56.43	18.52
	Trajectory	<u>76.02</u>	8.25	<u>14.64</u>	67.17	74.29	<u>16.49</u>
	Proceduralization	79.94	<u>9.76</u>	14.62	87.14	77.86	15.01
Claude-3.5-sonnet	No Memory	63.49	33.06	18.84	39.20	34.97	24.12
	Script	62.08	29.61	19.21	56.13	53.59	19.38
	Trajectory	<u>65.76</u>	29.61	<u>17.72</u>	<u>69.28</u>	<u>71.78</u>	<u>15.97</u>
	Proceduralization	65.46	<u>30.14</u>	15.29	82.50	74.72	15.79
Qwen2.5-72b	No Memory	56.57	7.34	18.32	44.91	41.25	21.38
	Script	58.59	7.34	18.53	<u>66.24</u>	61.88	17.13
	Trajectory	<u>63.41</u>	<u>12.66</u>	<u>18.12</u>	64.49	<u>69.57</u>	<u>16.40</u>
	Proceduralization	63.82	14.19	17.94	85.71	77.19	15.32

Table 1: Results on **Build Policy**. #CS, #HC denote Commonsense and Hard Constraint, respectively. \uparrow indicates the higher values are better, and \downarrow denotes the lower values are better. The best results among all methods with similar settings are **bolded**, and the second-best results are underlined.

of queries for retrieval.

- **AveFact**: We apply a large model to extract keywords from the task’s query, then computes the average similarity across matched keywords for retrieval.

During the retrieval process, we evaluate the similarity by calculating the cosine similarity between their corresponding vectors as shown in Table 2. Our experiments show that these different retrieval strategies produce varying results. Specifically, compared to random sampling, employing the query based and AveFact methods for precise retrieval significantly improves performance. The query-based approach benefits from capturing semantic contexts, enabling more accurate matches. The AveFact method, by extracting key features and averaging their similarities, effectively focuses on core task elements, leading to better retrieval efficacy. Overall, our findings suggest that incorporating semantic understanding and key feature extraction in retrieval strategies substantially enhances memory access accuracy and the effectiveness of downstream task performance.

4.3 Memory Update

While many prior efforts have focused on developing reusable procedural knowledge, enabling models to learn from prior experiences rather than solving each test task in isolation, most existing memory update methods remain quite rudimentary. Typically, they simply append newly acquired memories to the existing store—a so-called “merge”

Model	Policy	#CS \uparrow	#HC \uparrow	Steps \downarrow
GPT-4o	No Memory	71.93	12.88	17.84
	Random Sample	<u>74.59</u>	6.72	<u>15.12</u>
	Key=Query	73.38	<u>8.95</u>	15.44
	Key=AveFact	76.02	8.25	14.64
Claude-3.5-sonnet	No Memory	63.49	33.06	18.84
	Random Sample	63.99	<u>29.91</u>	17.93
	Key=Query	<u>64.93</u>	28.56	17.60
	Key=AveFact	65.76	29.61	<u>17.72</u>
Qwen2.5-72b	No Memory	56.57	7.34	18.32
	Random Sample	59.76	8.43	<u>18.31</u>
	Key=Query	<u>61.71</u>	<u>11.97</u>	18.54
	Key=AveFact	63.41	12.66	18.12

Table 2: Results on **Retrieve Policy** on TravelPlanner.

strategy. In this work, we explore several online memory-update mechanisms to identify which dynamic strategy delivers the best performance on our tasks. Beyond end-to-end evaluation metrics, we also analyze how both accuracy and efficiency evolve as the number of executed tasks increases, explicitly measuring the benefits conferred by our procedural memory.

To facilitate systematic comparison, we design several memory-update scenarios. The agent’s episodic memory is refreshed after every t test-set tasks, following the specific update strategies:

- **Vanilla Memory Update**: After every t tasks, all trajectories from these tasks are consolidated into procedural memories and directly appended to the memory bank.
- **Validation**: After every t tasks, the system performs a selective consolidation step: only the trajectories that terminated in successful

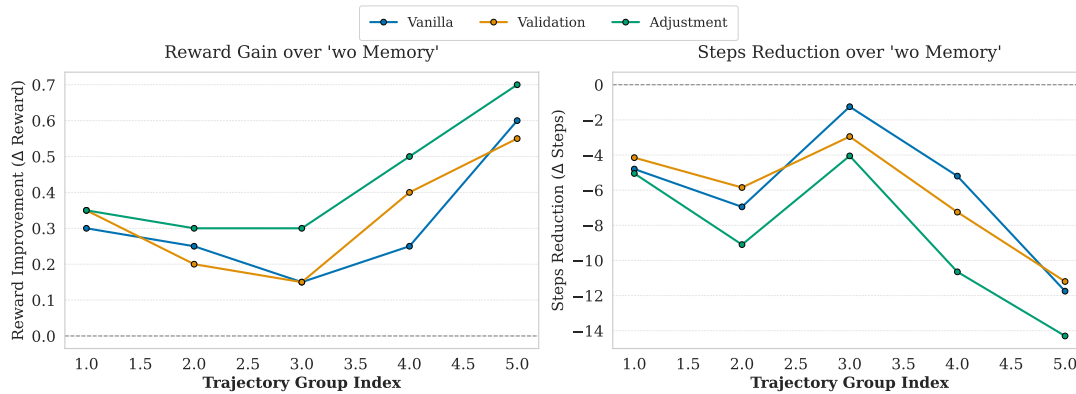


Figure 3: Reward gain and steps reduction vs. trajectory group index with **procedural memory**.

task completion are extracted, abstracted into compact, symbolic procedural memories. Trajectories that ended in failure, along with any redundant or noisy data, are discarded.

- **Adjustment:** When a retrieved procedural memory results in a failed execution, the erroneous trajectory is combined with the original memory and then revised in place, yielding an updated procedural memory.

As depicted in Figure 3, we systematically divided the tasks within our testbed into several distinct groups, with each group comprising a diverse set of individual tasks. Upon the completion of tasks within each group, we employed the previously described strategies to construct, store, and update the procedural memory. The experimental results reveal a clear trend: as we sequentially progress through more groups and iteratively refresh the memory, all strategies contribute to improved performance on subsequent tasks.

A closer comparison of different strategies exposes significant disparities in their effectiveness. Notably, the reflexion-based update mechanism stands out as the most effective approach. By the time the final group of tasks is reached, this method delivers a substantial advantage: it surpasses the second-best strategy by an impressive margin of +0.7 points and achieves a reduction of 14 steps. These improvements underscore the value of continually updating the memory, particularly when the update is guided by an error-correction mechanism embedded in the reflexion process.

4.4 Comparison

Figure 4 compares the performance and efficiency of different baselines using GPT-4o on ALFWorld.

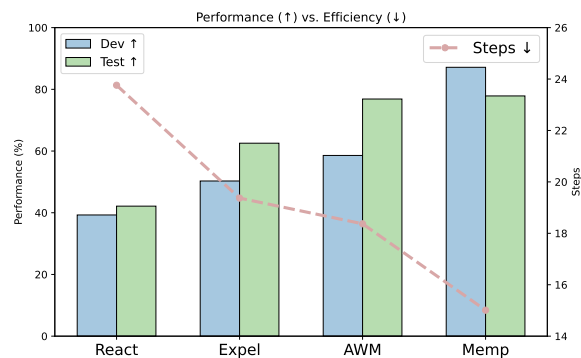


Figure 4: Performance of several baselines using GPT-4o on ALFWorld

React performs the worst, with low success rates and the highest number of steps, indicating inefficient interaction. Expel and AWM achieve progressively better performance, with AWM outperforming Expel on both splits. *Mem^p* delivers the best results, attaining the highest dev and test success rates while requiring the fewest steps. This shows that memory-augmented mechanisms significantly improve both task effectiveness and execution efficiency in long-horizon embodied agent tasks.

5 Analysis

Procedural Memory Boosts Accuracy and Cuts Trials. Figure 6 presents a case study demonstrating how Procedural Memory enhances both accuracy and efficiency. In the absence of Procedural Memory, facing a complex task that has not been performed before, there are usually two situations. In the first scenario, the model repeatedly attempts illegal or incorrect actions, causing the context to become increasingly complex and eventually exceeding the model’s understanding capacity. In the second scenario, after multiple attempts, the model

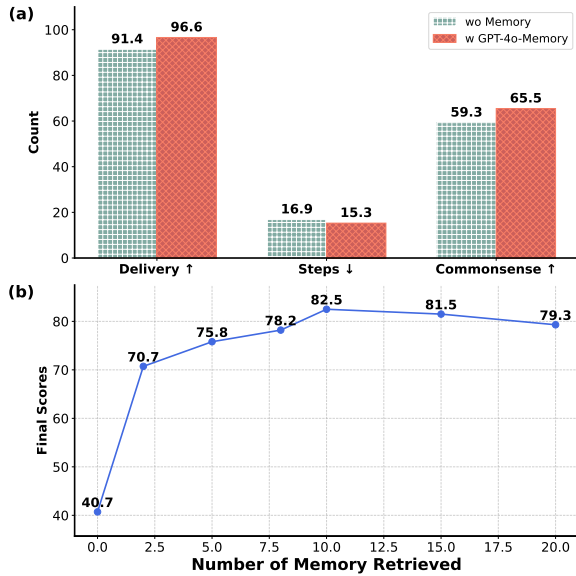


Figure 5: (a) Transfer result of GPT-4o’s procedural memory to Qwen2.5-14B-Instruct and its performance on TravelPlanner dataset. (b) The relationship between the quantity of procedural memory retrieved for GPT-4o’s performance on the ALFWorld dataset.

completes the task but at a cost significantly higher than the optimal path.

In contrast, once Procedural Memory is available for similar tasks, the model spends less time on trial and error. For example, in the egg heating problem, Procedural Memory can indicate the approximate location of the egg, saving the aimless search. During the heating process, it provides clear guidance, ensuring that the heating actions are performed consecutively and correctly, thereby allowing the task to be completed in fewer steps.

Procedural memory exhibits transferability from strong models to weaker ones. For a procedural memory constructed from a strong model in an offline memory library, we aim to verify whether this form of procedural memory can be effectively transferred to other models, or even weaker models. This exploration underscores the significance of memory transfer, as it could potentially enhance the adaptability and efficiency of various models by leveraging the knowledge and experience encapsulated within the strong model’s memory structure. As shown in Figure 5 (b), procedural memory generated by GPT-4o was employed by Qwen2.5-14B. On the Travel Plan benchmark, the 14 billion parameter model raised its task completion rate by 5% and cut the average number of steps by 1.6. Similar gains, both in success rate and trajectory

length, appeared on ALFWorld. These outcomes confirm that procedural knowledge from a stronger model can be distilled into a reusable memory bank and transferred to a smaller model with minimal overhead, giving that smaller model a clear boost in task solving ability.

Scaling Memory Retrieval Improves Agent Performance.

While our main experiment has already demonstrated that procedural memory improves an agent’s task accuracy and reduces the number of steps required, vector-based storage and retrieval confer an advantage over human procedural memory: they can be scaled both in total capacity and in the number of memories retrieved. To investigate whether an agent’s performance continues to rise as the procedural-memory store and the number of retrieved memories increase, we designed a set of follow-up experiments. As shown in Figure 5 (b), as the number of retrieved procedural memories increases, the agent’s performance also improves steadily, exhibiting an upward trend followed by a plateau. However, retrieving too many memories can lead to a decline in the agent’s performance. This is because excessive retrieval can affect the context length and also introduce less accurate procedural memories, which can interfere with the overall effectiveness.

6 Conclusion and Future Work

We introduce Mem^p , a task-agnostic framework that elevates procedural memory to a core optimization target in LLM-based agents. By systematically studying strategies for memory construction, retrieval, and updating, Mem^p enables agents to distill, reuse, and refine their own past experiences across diverse, long-horizon tasks and supports continual learning and robust generalization, marking a step toward self-improving, resilient agents.

In our experiments, Mem^p has achieved promising results in both construction and retrieval. Moving forward, we plan to enhance this work in several ways. Firstly, we will develop more diverse retrieval strategies. The current approach involves constructing different keys for vector-based retrieval. Secondly, in Mem^p , we currently rely on the standard reward signals provided by the benchmark. In such cases, using a LLM as a judge to assess task completion could be a viable solution. This would transform the agent’s lifecycle into a continuous loop of executing tasks, self-assessing completion, building memories, and then progress.

592 Limitations

593 *Mem^P* is still constrained by two key limitations:
594 1. retrieval is restricted to vector-similarity search
595 with manually crafted keys—other classic, poten-
596 tially more precise methods like BM25 have not
597 been incorporated; 2. the framework depends on
598 explicit, benchmark supplied reward signals, so it
599 cannot yet judge task success in real world settings
600 where such rewards are sparse or absent.

601 References

602 Anthropic. 2022. [Claude 3.5 sonnet system card](#).

603 Victor Barres, Honghua Dong, Soham Ray, Xujie Si,
604 and Karthik Narasimhan. 2025. τ^2 -bench: Evaluat-
605 ing conversational agents in a dual-control environ-
606 ment.

607 Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang,
608 Xingshan Zeng, Shuai Yu, Dexun Li, Shuai Wang,
609 Weinan Gan, Yuefeng Huang, and 1 others. 2025.
610 Acebench: Who wins the match point in tool usage?

611 Minghao Chen, Yihang Li, Yanting Yang, Shiyu Yu,
612 Binbin Lin, and Xiaofei He. 2024. [Automan-
613 ual: Constructing instruction manuals by llm agents
614 via interactive environmental learning](#). *Preprint*,
615 arXiv:2405.16247.

616 Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet
617 Singh, and Deshraj Yadav. 2025. Mem0: Building
618 production-ready ai agents with scalable long-term
619 memory.

620 Neal J Cohen and Larry R Squire. 1980. Preserved
621 learning and retention of pattern-analyzing skill in
622 amnesia: Dissociation of knowing how and knowing
623 that.

624 Guanting Dong, Yifei Chen, Xiaoxi Li, Jiajie Jin,
625 Hongjin Qian, Yutao Zhu, Hangyu Mao, Guorui
626 Zhou, Zhicheng Dou, and Ji-Rong Wen. 2025. Tool-
627 star: Empowering llm-brained multi-tool reasoner
628 via reinforcement learning.

629 Runnan Fang, Xiaobin Wang, Yuan Liang, Shuofei Qiao,
630 Jialong Wu, Zekun Xi, Ningyu Zhang, Yong Jiang,
631 Pengjun Xie, Fei Huang, and 1 others. 2025. Syn-
632 world: Virtual scenario synthesis for agentic action
633 knowledge refinement.

634 Erhu Feng, Wenbo Zhou, Zibin Liu, Le Chen, Yun-
635 peng Dong, Cheng Zhang, Yisheng Zhao, Dong Du,
636 Zhichao Hua, Yubin Xia, and 1 others. 2025. Get
637 experience from practice: Llm agents with record &
638 replay.

639 Prahlad Gupta and Neal J Cohen. 2002. Theoretical and
640 computational analysis of skill learning, repetition
641 priming, and procedural memory.

Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa
Safdari, Yutaka Matsuo, Douglas Eck, and Aleksan-
dra Faust. 2023. A real-world webagent with plan-
ning, long context understanding, and program syn-
thesis.

Yuki Hou, Haruki Tamoto, and Homei Miyashita. 2024.
" my agent understands me better": Integrating dy-
namic human-like memory recall and consolidation
in llm-based agents. In *Extended Abstracts of the
CHI Conference on Human Factors in Computing
Systems*, pages 1–7.

Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu,
Wenqi Shao, and Ping Luo. 2024a. Hiagent: Hier-
archical working memory management for solving
long-horizon agent tasks with large language model.

Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu,
Wenqi Shao, and Ping Luo. 2024b. Hiagent: Hier-
archical working memory management for solving
long-horizon agent tasks with large language model.

Tal Ifargan, Lukas Hafner, Maor Kern, Ori Alcalay,
and Roy Kishony. 2025. Autonomous llm-driven
research—from data to human-verifiable research
papers.

John E Laird. 2022. Introduction to the soar cognitive
architecture.

Wei Lan, Zhentao Tang, Mingyang Liu, Qingfeng Chen,
Wei Peng, Yiping Phoebe Chen, and Yi Pan. 2025.
The large language models on biomedical data analy-
sis: a survey.

Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii
Khizbullin, and Bernard Ghanem. 2023. Camel:
Communicative agents for" mind" exploration of
large language model society.

Zhiyu Li, Shichao Song, Chenyang Xi, Hanyu Wang,
Chen Tang, Simin Niu, Ding Chen, Jiawei Yang,
Chunyu Li, Qingchen Yu, and 1 others. 2025.
Memos: A memory os for ai system.

Bang Liu, Xinfeng Li, Jiayi Zhang, Jinlin Wang, Tanjin
He, Sirui Hong, Hongzhang Liu, Shaokun Zhang,
Kaitao Song, Kunlun Zhu, and 1 others. 2025a. Ad-
vances and challenges in foundation agents: From
brain-inspired intelligence to evolutionary, collabora-
tive, and safe systems.

Yitao Liu, Chenglei Si, Karthik Narasimhan, and
Shunyu Yao. 2025b. Contextual experience replay
for self-improvement of language agents.

Zexi Liu, Jingyi Chai, Xinyu Zhu, Shuo Tang, Rui Ye,
Bo Zhang, Lei Bai, and Siheng Chen. 2025c. MI-
agent: Reinforcing llm agents for autonomous ma-
chine learning engineering.

Fanbin Lu, Zhisheng Zhong, Shu Liu, Chi-Wing Fu,
and Jiaya Jia. 2025. Arpo: End-to-end policy opti-
mization for gui agents with experience replay.

695	Run Luo, Lu Wang, Wanwei He, and Xiaobo Xia. 2025.	Xiaoyu Tan, Bin Li, Xihe Qiu, Chao Qu, Wei Chu,	749
696	Gui-r1: A generalist r1-style vision-language action	Yinghui Xu, and Yuan Qi. 2025. Meta-agent-	750
697	model for gui agents.	workflow: Streamlining tool usage in llms through	751
698	Vasilios Mavroudis. 2024. Langchain v0. 3.	workflow construction, retrieval, and refinement .	752
699	OpenAI. 2022. Gpt-4 system card .	In <i>Companion Proceedings of the ACM on Web Con-</i>	753
700	OpenAI. 2025. Deep research system card .	ference 2025, WWW '25, page 458–467, New York,	754
701	Yixin Ou, Yujie Luo, Jingsheng Zheng, Lanning Wei,	NY, USA. Association for Computing Machinery.	755
702	Shuofei Qiao, Jintian Zhang, Da Zheng, Huajun	Xiangru Tang, Tianrui Qin, Tianhao Peng, Ziyang Zhou,	756
703	Chen, and Ningyu Zhang. 2025. Automind: Adap-	Daniel Shao, Tingting Du, Xinming Wei, Peng Xia,	757
704	tive knowledgeable agent for automated data science.	Fang Wu, He Zhu, Ge Zhang, Jiaheng Liu, Xingyao	758
705	Martin L Puterman. 1990. Markov decision processes.	Wang, Sirui Hong, Chenglin Wu, Hao Cheng, Chi	759
706	Shuofei Qiao, Runnan Fang, Zhisong Qiu, Xiaobin	Wang, and Wangchunshu Zhou. 2025. Agent kb:	760
707	Wang, Ningyu Zhang, Yong Jiang, Pengjun Xie, Fei	Leveraging cross-domain experience for agentic prob-	761
708	Huang, and Huajun Chen. 2025. Benchmarking agentic	lem solving . <i>Preprint</i> , arXiv:2507.06229.	762
709	workflow generation . In <i>The Thirteenth Inter-</i>	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Man-	763
710	national Conference on Learning Representations,	dlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and	764
711	<i>ICLR 2025, Singapore, April 24-28, 2025</i> . OpenRe-	Anima Anandkumar. 2023. Voyager: An open-ended	765
712	view.net.	embodied agent with large language models.	766
713	Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen,	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	767
714	Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang,	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	768
715	and Huajun Chen. 2023. Reasoning with language	Xu Chen, Yankai Lin, and 1 others. 2024a. A survey	769
716	model prompting: A survey . In <i>Proceedings of the</i>	on large language model based autonomous agents.	770
717	61st Annual Meeting of the Association for Comput-	Zhenhailong Wang, Haiyang Xu, Junyang Wang,	771
718	ational Linguistics (Volume 1: Long Papers), <i>ACL</i>	Xi Zhang, Ming Yan, Ji Zhang, Fei Huang, and Heng	772
719	2023, Toronto, Canada, July 9-14, 2023, pages 5368–	Ji. 2025. Mobile-agent-e: Self-evolving mobile assis-	773
720	5393. Association for Computational Linguistics.	tant for complex tasks.	774
721	Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang,	Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and	775
722	Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li,	Graham Neubig. 2024b. Agent workflow memory.	776
723	Yunxin Li, Shijue Huang, and 1 others. 2025. Ui-	Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin,	777
724	tars: Pioneering automated gui interaction with native	Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun	778
725	agents.	Xi, Gang Fu, Yong Jiang, and 1 others. 2025a. Web-	779
726	Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté,	dancer: Towards autonomous information seeking	780
727	Yonatan Bisk, Adam Trischler, and Matthew J.	agency.	781
728	Hausknecht. 2021. Alfworld: Aligning text and em-	Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang,	782
729	bedded environments for interactive learning . In <i>9th</i>	Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He,	783
730	International Conference on Learning Representa-	Deyu Zhou, Pengjun Xie, and Fei Huang. 2025b.	784
731	tions, <i>ICLR 2021, Virtual Event, Austria, May 3-7,</i>	WebWalker: Benchmarking LLMs in web traversal .	785
732	2021. OpenReview.net.	In <i>Proceedings of the 63rd Annual Meeting of the</i>	786
733	Hongjin Su, Ruoxi Sun, Jinsung Yoon, Pengcheng Yin,	Association for Computational Linguistics (Volume 1:	787
734	Tao Yu, and Sercan Ö Arık. 2025. Learn-by-interact:	Long Papers), pages 10290–10305, Vienna, Austria.	788
735	A data-centric framework for self-adaptive agents in	Association for Computational Linguistics.	789
736	realistic environments.	Xin Wu, Yuqi Bu, Yi Cai, and Tao Wang. 2024. Up-	790
737	Theodore Sumers, Shunyu Yao, Karthik Narasimhan,	dating large language models’ memories with time	791
738	and Thomas Griffiths. 2023a. Cognitive architectures	constraints.	792
739	for language agents.	x.ai. 2025. Grok 3 beta — the age of reasoning agents .	793
740	Theodore Sumers, Shunyu Yao, Karthik Narasimhan,	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen	794
741	and Thomas Griffiths. 2023b. Cognitive architectures	Ding, Boyang Hong, Ming Zhang, Junzhe Wang,	795
742	for language agents.	Senjie Jin, Enyu Zhou, and 1 others. 2025. The rise	796
743	Jingkai Sun, Qiang Zhang, Yiqun Duan, Xiaoyang Jiang,	and potential of large language model based agents:	797
744	Chong Cheng, and Renjing Xu. 2024. Prompt, plan,	A survey.	798
745	perform: Llm-based humanoid control via quantized	Menglin Xia, Victor Rühle, Saravan Rajmohan, and	799
746	imitation learning. In <i>2024 IEEE International Con-</i>	Reza Shokri. 2025. Minerva: A programmable mem-	800
747	ference on Robotics and Automation (ICRA), pages	ory test benchmark for language models.	801
748	16236–16242. IEEE.		

802	Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. 2024. Travelplanner: A benchmark for real-world planning with language agents . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	858
803		859
804		860
805		861
806		862
807		863
808	Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. 2025. A-mem: Agentic memory for llm agents.	864
809		
810		
811	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024a. Qwen2.5 technical report.	
812		
813		
814		
815	Hui Yang, Sifu Yue, and Yunzhong He. 2023. Auto-gpt for online decision making: Benchmarks and additional opinions.	
816		
817		
818	Yijun Yang, Tianyi Zhou, Kanxue Li, Dapeng Tao, Lu-song Li, Li Shen, Xiaodong He, Jing Jiang, and Yuhui Shi. 2024b. Embodied multi-modal agent trained by an llm from a parallel textworld. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 26275–26285.	
819		
820		
821		
822		
823		
824	Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik R Narasimhan. 2025. τ -bench: A benchmark for tool-agent-user interaction in real-world domains. In <i>The Thirteenth International Conference on Learning Representations</i> .	
825		
826		
827		
828		
829	Hongli Yu, Tinghong Chen, Jiangtao Feng, Jiangjie Chen, Weinan Dai, Qiyang Yu, Ya-Qin Zhang, Wei-Ying Ma, Jingjing Liu, Mingxuan Wang, and Hao Zhou. 2025. Memagent: Reshaping long-context llm with multi-conv rl-based memory agent . <i>Preprint</i> , arXiv:2507.02259.	
830		
831		
832		
833		
834		
835	Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. 2024. A survey on the memory mechanism of large language model based agents.	
836		
837		
838		
839	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. A survey of large language models.	
840		
841		
842		
843	Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. In <i>The Twelfth International Conference on Learning Representations</i> .	
844		
845		
846		
847		
848	Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. 2024a. Memorybank: Enhancing large language models with long-term memory. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19724–19731.	
849		
850		
851		
852		
853	Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. 2024b. Memorybank: Enhancing large language models with long-term memory. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19724–19731.	
854		
855		
856		
857		
	Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, Shiding Zhu, Jiyu Chen, Wentao Zhang, Xiangru Tang, Ningyu Zhang, Huajun Chen, Peng Cui, and Mrinmaya Sachan. 2023. Agents: An open-source framework for autonomous language agents . <i>Preprint</i> , arXiv:2309.07870.	858
		859
		860
		861
		862
		863
		864
	Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen, Shuai Wang, Xiaohua Xu, Ningyu Zhang, Huajun Chen, and Yuchen Eleanor Jiang. 2024. Symbolic learning enables self-evolving agents . <i>Preprint</i> , arXiv:2406.18532.	865
		866
		867
		868
		869
		870
	Zijian Zhou, Ao Qu, Zhaoxuan Wu, Sunghwan Kim, Alok Prakash, Daniela Rus, Jinhua Zhao, Bryan Kian Hsiang Low, and Paul Pu Liang. 2025. Mem1: Learning to synergize memory and reasoning for efficient long-horizon agents . <i>Preprint</i> , arXiv:2506.15841.	871
		872
		873
		874
		875
		876

877 **Appendix**

878 **A Use of AI Assistant**

879 We affirm that Large Language Models are em-
880 ployed solely as an assisted tool to refine wording
881 and sentence structure during our paper writing
882 process. Their use in the experiments is strictly for
883 scientific research purposes, and all such usage has
884 been explicitly documented in our Experimental
885 Settings and Reproducibility Statement. No other
886 reliance on LLMs is involved in this work.

887 **B Case Study**

888 Detailed case studies are provided in Figure 6.

889 **C Evaluation Details**

890 For ALFWorld dataset, task completion is evalu-
891 ated by the execution environment. After a task
892 is completed or the maximum number of execu-
893 tion steps is reached, the environment provides a
894 reward of 0 or 1 to indicate whether the task has
895 been successfully completed. For TravelPlanner,
896 we conduct experiments on the test set in a two-
897 stage mode. After multiple rounds of interaction
898 to obtain the travel trajectory and the final planner,
899 GPT-4o converts the travel plan into a specified
900 JSON format. The converted plan is then com-
901 pared with the gold standard to obtain scores for
902 both Common Sense and Hard Constraint.

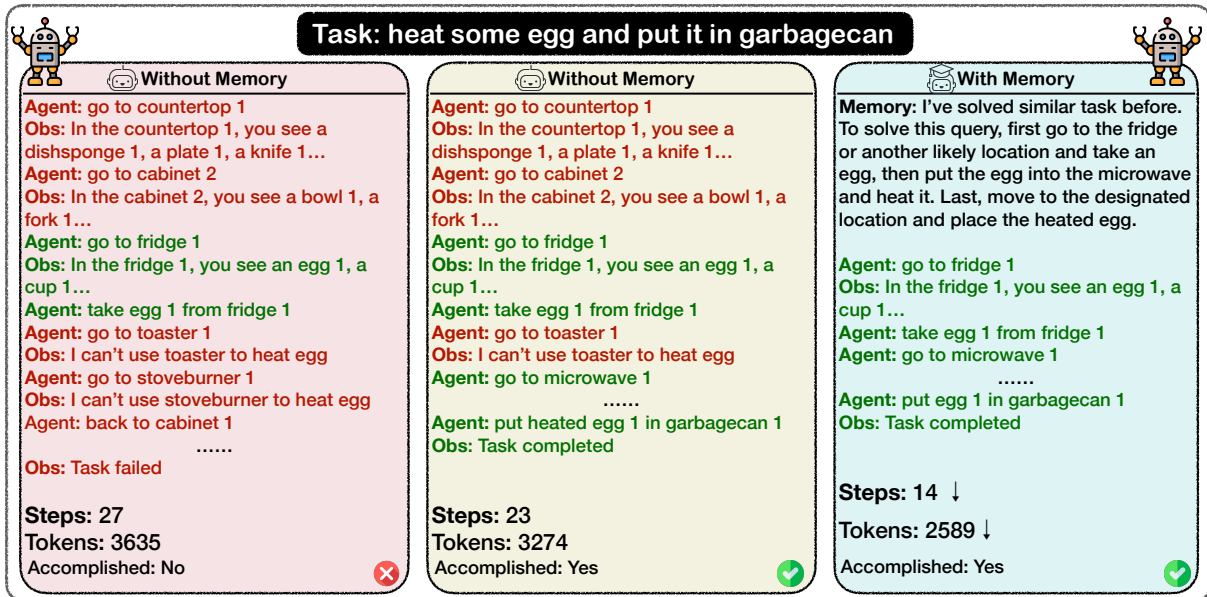


Figure 6: Compare trajectories with and without procedural memory, shortens the process by 9 steps and saves 685 tokens.