# MEM1: Learning to Synergize Memory and Reasoning for Efficient Long-Horizon Agents

**Zijian Zhou**[*12]    **Ao Qu**[*13]    **Zhaoxuan Wu**[1]    **Sunghwan Kim**[4]
**Alok Prakash**[1]    **Daniela Rus**[13]    **Jinhua Zhao**[13]    **Bryan Kian Hsiang Low**[13]
**Paul Pu Liang**[3]
[1]Singapore-MIT Alliance for Research and Technology Centre
[2]National University of Singapore
[3]MIT
[4]Yonsei University

## Abstract

Modern language agents must operate over long-horizon, multi-turn interactions, where they retrieve external information, adapt to observations, and answer interdependent queries. Yet, most LLM systems rely on full-context prompting, appending all past turns regardless of their relevance. This leads to unbounded memory growth, increased computational costs, and degraded reasoning performance on out-of-distribution input lengths. We introduce **MEM1**, an end-to-end reinforcement learning framework that enables agents to operate with **constant memory** across long multi-turn tasks. At each turn, MEM1 updates a **compact shared internal state** that jointly supports memory consolidation and reasoning. This state integrates prior memory with new observations from the environment while strategically discarding irrelevant or redundant information. To support training in more realistic and compositional settings, we propose a simple yet effective and scalable approach to constructing multi-turn environments by composing existing datasets into arbitrarily complex task sequences. Experiments across three domains, including internal retrieval QA, open-domain web QA, and multi-turn web shopping, show that MEM1-7B improves performance by $3.5\times$ while reducing memory usage by $3.7\times$ compared to Qwen2.5-14B-Instruct on a 16-**objective** multi-hop QA task, and **generalizes beyond the training horizon**. Our results demonstrate the promise of reasoning-driven memory consolidation as a scalable alternative to existing solutions for training long-horizon interactive agents, where both efficiency and performance are optimized. Code can be found at https://github.com/MIT-MI/MEM1.

## 1 Introduction

Large language models (LLMs) have shown remarkable performance in single-turn tasks such as question answering, summarization, and code generation [7, 51, 3]. However, emerging real-world applications increasingly operate over multiple turns—searching documents, interacting with environments [70], and making decisions based on evolving external information [53]. Examples include research agents such as OpenAI and Gemini Deep Research [37, 18] that automate complex tasks by iteratively gathering information, and web-navigation agents such as OpenManus [50] and BrowserUse [34], which must complete goals across dozens of interactive turns.

Unlike traditional tasks where the input is static or self-contained, long-horizon settings often involve answering a sequence of related questions, requiring the agent to continuously retrieve new

---

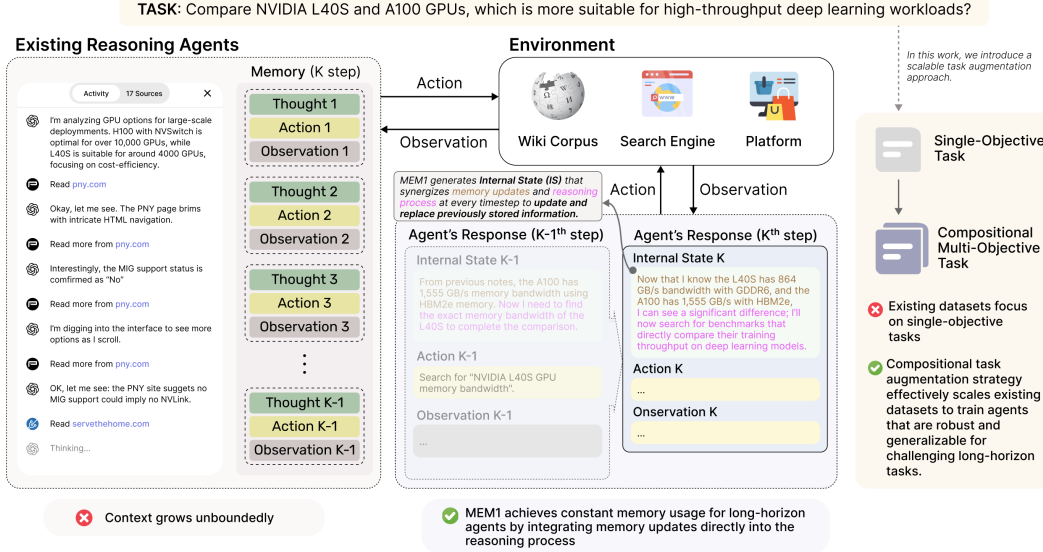* Equal contribution. Correspondence: zhou_zijian@u.nus.edu, qua@mit.edu

Figure 1: Comparison of memory management between MEM1 and existing reasoning agents. While existing agents for long-horizon tasks [24, 64, 69] continuously append thoughts (typically enclosed in <think></think>), actions, and observations, resulting in an ever-growing context, our MEM1 agent learns to keep updating an internal state (enclosed in <IS></IS>) that blends thought and memory, discarding the contents from previous steps to achieve constant memory usage during the task. On the other hand, while existing environments and datasets focus on single-objective tasks, our task augmentation method effectively scales up these tasks to enable long-horizon agent training.

information, revise beliefs, and adapt to evolving contexts over time. For instance, consider a research assistant tasked with "What's the evidence for X?". Subsequent queries like "Who published it?" require further information retrieval, while "Is the source credible?" calls for self-reflection and assessment. Each query builds on the previously collected and accumulated information. Similarly, a shopping assistant may be first asked "Which product is cheapest?", then "What are its reviews?", and "Is it compatible with my device?". These interactions span multiple turns, featuring evolving contexts and compound reasoning.

In systems designed for long-horizon settings, a common approach is to append all past observations, actions, and thoughts to the context at every turn [55, 61]. This forces the model to operate with an unboundedly growing context, which introduces three key challenges. **(1) Growing inference cost and memory usage**. Transformer-based LLMs typically incur $O(N^2)$ compute cost (or $O(N)$ with Key-Value caching) and $O(N)$ memory usage as the context length $N$ increases [52]. Consequently, deploying these models requires reserving large GPU memory on modern inference frameworks to accommodate the growing context [27, 68], often leading to significant wastage of computing resources. **(2) Generalization limits beyond the training horizon**. Ongoing conversations with context length exceeding that in the training data become out-of-distribution for the model. The model struggles to manage and reason over such unfamiliar long-horizon inputs [63]. **(3) Overloaded and inefficient context**. The accumulation of irrelevant or redundant content dilutes the model's attention. This reduces its ability to reason effectively, even when relevant information is technically still present within the prompt [2, 30, 56].

Recent progress in long-context modeling largely targets static inputs (*e.g.,* long documents) and does not address multi-turn interaction with external environments [6, 19]. Some other approaches introduce external memory modules (*e.g.,* summarizers or retrievers) [63, 29, 13, 57], but these are typically trained separately and cannot be optimized end-to-end with the agent's policy. This also introduces additional engineering overhead, as engineers must manage and integrate two separate models. Meanwhile, existing works on tool-using agent systems trained with reinforcement learning leave memory management unsolved, letting the prompt length grow unboundedly [24, 69]. A natural question is raised: *Can a language model learn to consolidate its memory as part of its reasoning process* so that it retains only what is essential for solving the task?

Motivated by this question, we present MEM1: **M**emory-**E**fficient **M**echanism via learning **1**-step integrated reasoning and consolidation—a method for training LLM agents that maintain constant memory usage across arbitrarily long horizons. As illustrated in Fig. 1, at each turn, the model updates a consolidated state composed of prior memory and newly obtained information. This consolidated state becomes the agent's only retained memory, allowing all external tool outputs to be discarded after use, which prevents prompt expansion altogether as illustrated by Fig. 2. A key insight of our method is that inference-time reasoning [55, 15, 33, 65] serves a dual function: it not only provides deeper insight into the current query but also acts as a form of "working memory" [4], extracting key components from gathered information to build an evolving understanding. By unifying reasoning and memory consolidation, MEM1 offers an elegant solution where the agent learns to reason and remember within a shared representational space without requiring additional modules or architectural changes.

We train this behavior end-to-end with reinforcement learning (RL) [48, 71], optimizing for task success via verifiable rewards [43]. Although not explicitly optimized for memory efficiency through reward signals, the agent learns to manage memory as part of its policy, resulting in near-constant memory usage across long horizons. Additionally, we notice that current training and evaluation environments predominantly focus on single-objective tasks [26, 59, 38], limiting their ability to fully prepare agents for realistic, long-horizon scenarios that inherently involve multiple sequential objectives. To address this challenge, we introduce a scalable task augmentation approach, transforming existing single-objective QA datasets into complex multi-objective tasks through compositions of $N$ multi-hop questions. This approach enables us to repurpose standard benchmarks in our community to more effectively train and evaluate agents on long-horizon reasoning, an increasingly important capability in real-world applications.



Figure 2: A conceptual comparison of context length between the MEM1 agent and existing reasoning agents when handling long-horizon tasks. Our agent learns to discard the previous context (except for the prompt and initial query) immediately after generating a new internal state and action, resulting in near-constant memory usage.

To evaluate our method comprehensively, we employ diverse multi-turn environments, including internal retrieval-augmented QA [26, 59], open-domain Web QA [69], and complex multi-turn agent shopping scenarios in WebShop [60]. Across these scenarios, MEM1 consistently rivals the performance of leading baselines while delivering efficiency gains up to $3.5\times$ in memory usage. Moreover, agents trained on our augmented 2-objective compositions generalize effectively to tasks involving up to 16-objective compositions. Notably, at the 16-objective level, our MEM1 achieves superior accuracy compared to all baseline methods, along with $1.27\times$ lower peak memory usage and $1.78\times$ faster inference compared to the respective best uncollapsed baseline.

## 2  Related Work

**LLM agents in multi-turn environment.**  LLM-based agents have evolved from handling single-turn queries to serving as autonomous agents capable of multi-turn interactions such as web navigation [60, 70] and complex research [69]. To enable such capabilities, Yao et al. [61] introduced the ReAct (*i.e.,* Reason + Act) framework, which enhances LLMs' ability to interact with external environments by interleaving reasoning and action. Building on this reasoning-acting prompting paradigm, subsequent works have explored ways to improve agent performance through natural language feedback, enabling iterative refinement [45, 32]. Recently, inference-time scaling has emerged as a promising direction for enabling complex reasoning, with prior research incorporating evaluators (*e.g.,* verifier, reward model) [46, 31] or world models [9]. In addition, there are two major
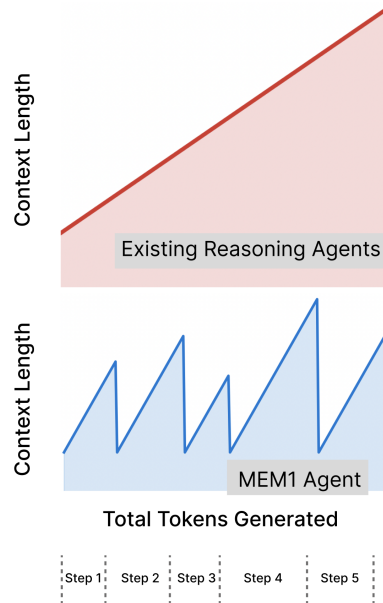
lines of training approaches: (1) behavior cloning (BC), which involves imitating expert trajectories to guide agent behavior by supervised fine-tuning (SFT) [62, 16, 12], and (2) reinforcement learning (RL), which optimizes agent policies by incentivizing desirable outcomes through rewards [47, 5, 39]. These methods aim to align the agents' behaviors with task objectives, enabling more robust and generalizable performance.

**Memory management for LLM agents.** A widely adopted approach to memory management in LLM-based agent systems involves appending all prior information, such as observations, intermediate thoughts, and actions, into the prompt at each interaction turn [61]. While this method is straightforward and effective when the number of interactions required is small, it results in unbounded context growth, leading to linearly scaled inference memory. Moreover, long contexts often contain irrelevant or redundant information, which impairs the model's reasoning capabilities [2, 30, 56]. To mitigate these issues, recent studies have proposed external memory frameworks, including retrieval-augmented generation and summarization modules [63, 29, 13, 57]. However, these methods are typically trained or used independently of the agent's policy, creating a disconnect between memory and the reasoning process. In addition, managing and integrating such modules often incurs extra computational overhead and system complexity. Despite these advancements, many RL approaches for training LLM agents still rely on accumulating the full interaction history as memory [24, 69, 39], leaving memory management during training an underexplored area. In this work, we seek to bridge this gap by tightly integrating memory with the agent's reasoning process, thereby enabling more efficient and context-aware decision-making.

## 3 MEM1

Complex reasoning tasks often require an iterative process of information gathering and synthesis, as seen in applications such as "deep research" [37, 23] and web-based agents [35, 20]. Recent advances in agent design involve interaction loops that interleave chain-of-thought reasoning [55, 15], environment interaction, and real-world feedback collection. To explicitly capture these core elements, we annotate each component using XML-style tags: `<IS>` for internal state (reasoning), `<query>` for environment queries, `<answer>` for the agent's responses, and `<info>` for external observations or tool outputs. MEM1 adopts a **learned** approach to iterative state updating and consolidation, ensuring that only the most recent set of `<IS>`, `<query>`, `<answer>`, and `<info>` elements is retained in the prompt at any given time. This design maintains a bounded and semantically relevant context, promoting efficient and coherent multi-step reasoning.

### 3.1 Memory as Part of Reasoning

To achieve a constant memory, MEM1 is particularly trained to iteratively refine its understanding by processing new information in conjunction with a consolidation of its prior state. At each turn $t$, the agent produces a new `<IS_t>` element, which summarizes past information and reasons about subsequent actions. Following this, the agent generates an action—either a `<query_t>` to interact with the environment, or an `<answer_t>` if a direct response is warranted. If the agent issues a `<query_t>`, the corresponding feedback from the environment is appended as `<info_t>`. At the next turn, $t + 1$, the agent consolidates the tuple $\left(\texttt{<IS\_t>}, \texttt{<query\_t>}, \texttt{<info\_t>}\right)$ into a new `<IS_(t+1)>`, which serves as the basis for further interactions. After each turn, all tags from the previous turn $t$ are pruned from the context, effectively compressing memory and preventing prompt bloat. Fig. 3 (bottom left) illustrates the evolution of the model's context over time. At any given turn, the agent retains at most two `<IS>` elements, two `<query>` elements, and one `<info>` element, ensuring bounded and efficient memory usage. The detailed rollout algorithm is in Alg. 1 of App. A.5.

RL offers a powerful mechanism for shaping agent behavior through reward signals [49]. In MEM1, we leverage this framework to incentivize effective state consolidation by designing environments in which the agent is rewarded only when it strategically retains and integrates useful information. Specifically, we construct tasks that require numerous interactions with the environment to arrive at a correct answer (see Sec. 3.3). Success depends on the agent's ability to rely on information collected along the inference path. At each turn, we prune the agent's context to retain only the most recent `<IS>`, forcing the agent to perform memory consolidation as part of its reasoning process. Without access to full historical context, the agent must learn to preserve and update relevant knowledge internally in order to reap the reward. This learning procedure mirrors how humans
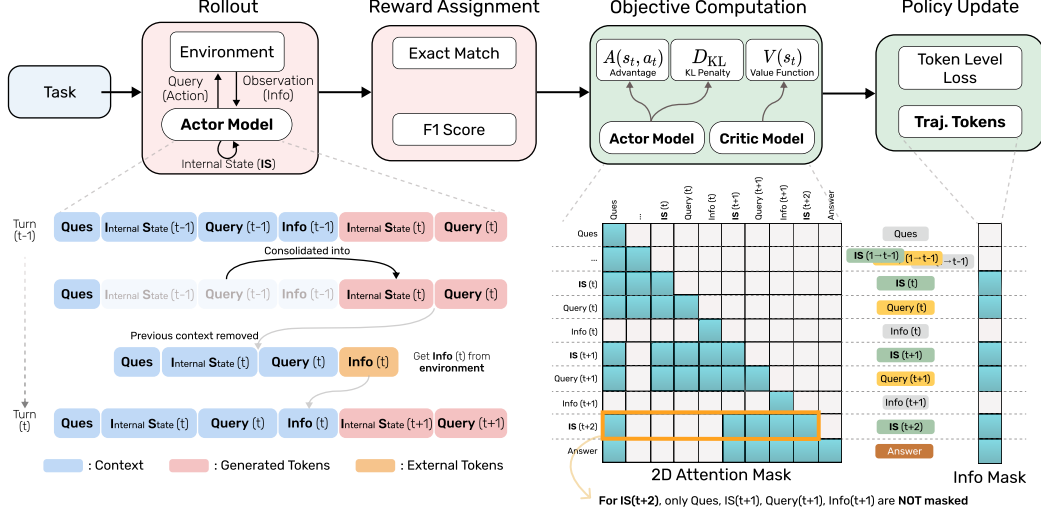
4

Figure 3: (Top): the RL pipeline used to train MEM1. (Bottom left): The evolution of context in MEM1–old `<IS>`, `<query>`, `<info>` are cleared as new states enter the context. The mechanism is used in the rollout. (Bottom right): the 2D attention mask used during the objective computation stage. The mask is applied during the forward pass to compute action log-probabilities for the actor model and state value estimates for the critic model. During the policy update stage, the information mask is then applied to the full trajectory, masking out tokens that were not generated by the model itself.

cultivate memorization skills through structured tasks such as Sudoku or crosswords [1], where success hinges on selectively attending to key information and building upon it. Over time, such tasks help individuals develop cognitive strategies that jointly support efficient memorization and reasoning, similar to our RL method for training MEM1.

## 3.2 Masked Trajectory for Policy Optimization

In the previous section, we detailed the rollout process of our RL pipeline. However, unlike conventional multi-turn agents that preserve a static context during generation, MEM1 introduces a unique challenge: its mechanism continuously updates the context at each turn by consolidating prior memory and pruning irrelevant tokens. This dynamic context update disrupts the continuity of the token generation trajectory, complicating the estimation of token-wise advantages in policy optimization algorithms such as PPO and Reinforce++ [41, 22], where trajectories are typically assumed to be linear.

To address this, we introduce a masked trajectory approach that reconstructs a logically coherent full trajectory by stitching together multiple interaction turns with evolving contexts. This unified trajectory mimics a standard multi-turn rollout and comprises a sequence of tuples $\tau_t = (\texttt{<IS\_t>}, \texttt{<query\_t>}, \texttt{<info\_t>})$ for $t \in [1, T-1]$, where $T$ denotes the total number of interaction turns. The $T$-th turn outputs the final answer $\tau_T = (\texttt{<IS\_t>}, \texttt{<answer\_t>})$. The full trajectory encodes all information needed for accurate policy learning while respecting MEM1's memory consolidation at each turn. Fig. 3 (bottom left) demonstrates the evolution of the agent's context.

To ensure that policy gradients are correctly computed under this consolidated memory regime, we apply a **two-dimensional attention mask** [40] across the full trajectory. This mask restricts each token's attention to only the tokens retained in memory at the time that token was generated. Specifically, for a token position $k$, we mask out all prior tokens that are not part of the consolidated memory corresponding to the context at that turn. This masking strategy enables accurate computation of the policy objective: letting $s_k$ denote the masked input state and $a_k$ the predicted token, the log-probability ratio $\rho_k(\theta) = \frac{\pi_\theta(a_k|s_k)}{\pi_{\theta_{\text{old}}}(a_k|s_k)}$ remains valid and tractable, in turn ensuring that the advantage, KL penalty, and value estimation are correct. A visualization is in Fig. 3 (bottom right). Furthermore, following [24], we incorporate a one-dimensional attention mask over retrieved external information

5

during policy updates for both the actor and critic networks. This ensures that gradient updates are localized to only tokens generated by the agent. Fig. 3 (bottom right) shows the masking mechanism that enables stable and accurate policy optimization under MEM1's memory-constrained execution.

## 3.3 Multi-Objective Task Design

Although our proposed method is designed to address the critical challenges of agentic multi-turn interaction with the external world, there are limited publicly available datasets that support training for such long-horizon interactive processes. Existing benchmarks, such as HotpotQA [59], Bamboogle [38], and 2wiki [21], are often cited as multi-hop benchmarks, yet they typically involve only two information-seeking steps. Moreover, these datasets are not explicitly structured to support long-horizon interactions that necessitate the agent to manage the memory state.

To bridge this gap, we introduce a novel task—multi-objective question & answering (QA)—that extends the number of reasoning steps required to solve a problem. Building on existing multi-turn datasets such as HotpotQA and Natural Question [59, 26], we interleave multiple questions from the original QA corpus and construct a single composite query that requires answering all constituent sub-questions, shown in Prompt 1 of App. A.3. This formulation compels the agent to perform multiple search queries, each targeting a distinct sub-question or sub-objective, and then integrate the retrieved answers to form a comprehensive final response. Compared to the original tasks, our synthesized multi-objective, multi-hop setting significantly increases the number of search and reasoning turns required, leading to more complex, memory-intensive interactions.

## 4 Experiments & Results

We empirically demonstrate the effectiveness of our approach in training the MEM1 agent to perform multi-turn tasks while preserving a near-constant-sized memory state. We evaluate MEM1 against several baselines using a comprehensive set of metrics categorized into *accuracy* (*e.g.,* Exact Match, F1 score, Environment Reward) and *efficiency* (e.g., Peak Token Usage, Dependency Length, Inference Time). All MEM1 variants are fine-tuned from the Qwen2.5-7B Base model [58]. We use PPO [41] as the RL algorithm as it computes token-level advantages, bringing stability to the training process. While we also experimented with instruction-tuned and supervised fine-tuned models using curated high-quality trajectories, reinforcement learning from the base model consistently yielded the best performance and generalization.

Our experiments are conducted in two standard environments, each reflecting real-world scenarios that require multi-turn agent interactions. The first environment is question answering with retrieval-augmented generation (RAG) [26, 59], where the agent must answer queries by retrieving relevant information from an external knowledge store (either a database or an online search engine). We trained on RAG with a local database (*i.e.,* Wikipedia Corpus) and evaluated on tasks involving open web browsing. For QA, following Sec. 3.3, we construct multi-objective tasks and tested the model performance on tasks with more questions than seen in the training. The second environment is WebShop navigation [60], where the agent assists users in online shopping by browsing a website and selecting items based on natural language descriptions. This task requires the agent to iteratively read page content and make navigation decisions, following protocols similar to those in WebGPT [35].

### 4.1 Implementation Details

**Datasets and evaluation metrics.** We train two versions of MEM1 agent for both long-horizon QA and web navigation. For long-horizon QA, we augment the multi-hop QA dataset from [24] that mixes data from both HotpotQA [59] and Natural Question [26] to form a 2-objective composite task.

For the web agent, we use the WebShop environment [60], which also produces a reward during training [64]. For all datasets, the train-test split follows the original papers. During RL training, we employ the exact match (EM) metric for QA tasks (details in App. A.4.1) and the environment reward for WebShop [60, 64]. To evaluate the effectiveness of various approaches, we measure the EM and F1 score for QA tasks and final reward for the WebShop environment [60, 64]. To evaluate the efficiency, we consider the peak token usage, average dependency, and average inference time. The test datasets are obtained from the original papers which consist of out-of-distribution data. The former two metrics measure the memory efficiency, while the latter measures the time efficiency.

The detailed definitions of the metrics are in App. A.4.1. The prompt and format can be found in App. A.3.

**Baselines.** To evaluate the accuracy and efficiency of MEM1, we compare it against a diverse set of baselines designed to either enhance task performance or manage context effectively. For the QA environment, we benchmark accuracy against Search-R1 [24], DeepResearcher [69], and a larger-scale model, Qwen2.5-14B-Instruct [58]. Details about Search-R1 and DeepResearcher can be found in App. A.4.2. For the WebShop environment, we compare against Agent-FLAN [11], Agent-R [64], and AgentLM [66]. To assess efficiency, we consider two context compression baselines using models of the same parameter size as MEM1. First, we apply MEM1's agentic *truncation* prompt template and rollout to a standard instruct model, isolating the benefits of prompt and rollout design alone. Second, we evaluate A-MEM [57], which augments an Instruct model with a vector database for memory retrieval, capturing the effect of external memory modules in agentic systems. We additionally train a supervised fine-tuned (SFT) model using trajectories curated from GPT-4o [36] based on MEM1's rollout and compare it with the RL-trained agent.

**Meta info injection.** In our agentic pipeline, the agent's context is programmatically truncated at each turn—immediately after it generates a search query or an answer—following the procedure outlined in Sec. 3. As past context is truncated, the agent may have difficulty determining when to terminate. To address this, we prepend a hint [HINT: YOU HAVE {turns_left} TURNS LEFT] at the beginning of each `<info>` tag to remind the agent of its remaining turns budget. For all experiments, we set the maximally allowed turns to 6 for 1-objective to 4-objective tasks and 20 for more difficult tasks to avoid excessively long trajectories.

## 4.2 MEM1 on Multi-Objective Multi-Hop Tasks

One key advantage of MEM1 agents lies in their efficient management of long-horizon interactions with the environment. To demonstrate this, we train our MEM1 agent with a 2-objective augmentation of the QA dataset, and subsequently test it against other models, using held-out multi-objective test datasets similarly augmented from the original test datasets. As elaborated in Sec. 3.3, these multi-objective tasks require a significantly larger number of turns of environment interactions to complete, hence serving as better benchmarks for memory management. As shown in Tab. 1, when evaluated on 2-objective datasets, MEM1 achieves better performance (in terms of EM and F1 scores) than other 7B counterparts, while incurring significantly lower peak token usage and achieving faster inference time.

The advantage of MEM1 becomes even more evident in tasks requiring longer-horizon interactive processes. To highlight such scalability of MEM1, we further compare the models on $3, 4, 6, 8,$ and 16-objective tasks in Fig. 4 and Tab. 1. Fig. 4 illustrates the scaling trends of task performance (measured by EM count) and memory efficiency (measured by Peak Token Usage) for MEM1 relative to other models and memory management baselines. As the number of objectives increases, the Peak Token Usage of all other methods and models scales nearly linearly. In contrast, MEM1 maintains an almost constant peak token count with only a slight increase, as also shown in Tab. 1.

Notably, while MEM1 initially underperforms Qwen2.5-14B-Instruct, its performance gradually catches up as the number of objectives increases, eventually surpassing the 14B model, which has double the parameter count. MEM1 also demonstrates remarkable efficiency. In the 16-objective task, it requires only $27.1\%$ of the peak tokens and $29.3\%$ of the total inference time compared to Qwen2.5-14B-Instruct. This efficiency translates to significantly reduced GPU memory requirements and overall computing resource demands.

## 4.3 MEM1 on Single-Objective Multi-Hop Tasks

While MEM1 is designed to train agents for very long-horizon tasks, our training method also delivers improved capability with existing multi-hop tasks while achieving much greater efficiency at the same time, all without being explicitly trained on the single-objective versions of these tasks. Note that single-objective tasks also require multiple turns of interaction to produce the desired output.

**Long-horizon web navigation in WebShop.** Beyond QA tasks, we further evaluate the effectiveness of MEM1 in managing long-horizon interactions in the form of web navigation. We show
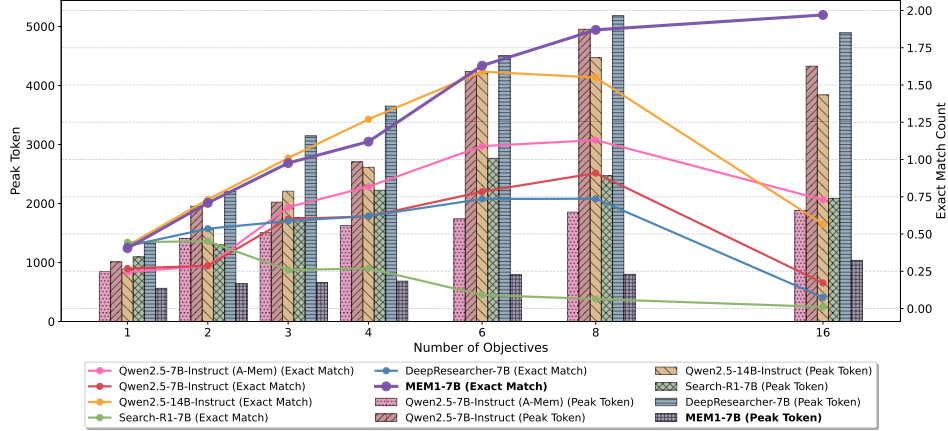
Figure 4: Performance and efficiency scaling of MEM1 (trained on 2-objective QA) with the number of objectives in multi-objective tasks. MEM1 outperforms the other models and baselines while having an almost constant scaling in memory usage. Note that at 16-objective, the context of baseline models does not increase anymore since their model performance has degraded (some collapsed).

Table 1: Comparison of models on multi-objective multi-hop QA tasks. Arrows indicate the desired directions. Numbers in red indicate collapsed model behavior (extremely low performance). (truncate) means using MEM1's prompt and rollout pipeline. (A-MEM) means using MEM1's prompt and rollout pipeline with A-Mem's external memory module [57]. MEM1-QA means MEM1 trained on 2-objective QA task.

| Model | 2-Objective | | | | 8-Objective | | | | 16-Objective | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | EM ↑ | F1 ↑ | Peak ($\times 10^2$) ↓ | Time (s) ↓ | EM ↑ | F1 ↑ | Peak ($\times 10^2$) ↓ | Time (s) ↓ | EM ↑ | F1 ↑ | Peak ($\times 10^2$) ↓ | Time (s) ↓ |
| Qwen2.5-14B-Inst | **0.732** | **0.902** | 15.6±0.19 | 5.49 ± 0.16 | 1.55 | 1.87 | 44.7 ± 0.37 | 16.2 ± 0.27 | 0.567 | 0.703 | 38.4±0.71 | 29.7±0.75 |
| Qwen2.5-7B-Inst | 0.268 | 0.366 | 19.6±0.33 | 4.60±0.08 | 0.87 | 1.10 | 49.5±0.40 | 13.9±0.18 | 0.165 | 0.213 | 43.3±0.62 | 15.5±0.23 |
| Qwen2.5-7B-Inst (A-MEM) | 0.286 | 0.371 | 14.1±0.10 | 24.6±0.51 | 1.13 | 1.43 | 18.6±0.10 | 53.7±1.26 | 0.730 | 0.961 | 18.8±0.14 | 91.2±2.44 |
| Qwen2.5-7B-Inst (truncate) | 0.262 | 0.336 | 8.28±0.06 | 5.89±0.16 | 0.97 | 1.23 | 11.8±0.10 | 11.9±0.20 | 0.396 | 0.497 | 13.3±0.16 | 22.1±0.60 |
| Search-R1 | 0.452 | 0.531 | 13.0±0.08 | 4.09 ± 0.23 | 0.064 | 0.08 | 24.7 ± 0.19 | 4.25±0.16 | 0.009 | 0.011 | 20.9±0.03 | 4.75±0.18 |
| DeepResearcher | 0.536 | 0.650 | 22.0±0.43 | 4.01±0.07 | 0.73 | 0.90 | 51.8±0.35 | 11.3±0.14 | 0.071 | 0.106 | 48.9±0.66 | 15.8±0.19 |
| **MEM1-QA** | 0.709 | 0.838 | **6.40±0.02** | 6.49 ± 0.07 | 1.87 | 2.31 | **8.01±0.06** | 8.68±0.12 | 1.97 | 2.39 | **10.4±0.09** | 8.70±0.12 |

the experimental results in Tab. 2. Trained in the WebShop environment (see App. A.6), MEM1 outperforms other agent training baselines, including Agent-Flan, Agent-R, and AgentLM when utilizing models of similar size. Furthermore, MEM1 achieves remarkable efficiency improvements compared to the best baseline method, AgentLM, featuring a $2.8\times$ improvement in Peak Token Usage, a $1.9\times$ improvement in Dependency, and a $1.5\times$ improvement in Inference Time. MEM1 even surpasses AgentLM-13B, a model with twice the parameter count of our trained model. Additionally, our results indicate that using MEM1 is significantly better than OpenAI's GPT-4o on the WebShop tasks, even when the truncation prompt templates or A-MEM techniques are applied to GPT-4o.

**Single-objective QA in Wikipedia.** Tab. 3 presents the accuracy and efficiency metrics for evaluations on single-objective QA tasks on Wikipedia [24], where the agent can make retrieval requests from the Wikipedia datastore via RAG. The MEM1 used in this evaluation is the same as the one detailed in Sec. 4.2, which is trained solely on a 2-objective task. Overall, MEM1 demonstrates superior efficiency across all three evaluated efficiency metrics, while simultaneously achieving the highest EM score and an F1 score comparable to that of Qwen2.5-14B-Instruct. This improvement in efficiency is attributed to the MEM1 agent's ability to consolidate memory from previous interactions into a compact internal state, which reduces the number of tokens used in the context. We also observe that SFT significantly underperforms RL, highlighting the necessity for RL-based training.

**Zero-shot transfer to Online Web-QA.** To validate the transferability and generalizability of the trained MEM1 agent, we perform a zero-shot transfer to an online web-QA environment, which is unseen by the agent. In this environment, agents conduct web searches through an API service that returns results including titles, snippets, and URLs. As shown in Tab. 3, MEM1 consistently exhibited improved efficiency alongside comparable effectiveness in this unseen setting via zero-shot transfer.

Table 2: The experimental results for WebShop. For a fair comparison, we do not report GPT's inference time. For Agent-R, scores are taken from the original paper, as the model is closed source. MEM1-WebShop means MEM1 trained on WebShop environment.

| Model | Avg Final Reward ↑ | Peak Token ($\times 10^3$) ↓ | Dependency ($\times 10^6$) ↓ | Inference Time Per Traj (s) ↓ |
|---|---|---|---|---|
| GPT-4o | 25.48 | $5.30 \pm 1.23$ | $3.99 \pm 1.16$ | *N/A* |
| GPT-4o (truncate) | 13.82 | $0.99 \pm 0.99$ | $0.81 \pm 0.23$ | *N/A* |
| GPT-4o (A-MEM) | 24.50 | $1.84 \pm 0.06$ | $0.31 \pm 0.11$ | *N/A* |
| Qwen2.5-7B-Instruct | 18.42 | $5.64 \pm 1.34$ | $3.38 \pm 0.89$ | $12.31 \pm 1.82$ |
| Qwen2.5-14B-Instruct | 12.34 | $5.44 \pm 0.92$ | $3.30 \pm 0.61$ | $18.17 \pm 2.32$ |
| Agent-FLAN-7B | 40.35 | $3.37 \pm 1.12$ | $2.18 \pm 1.62$ | $9.95 \pm 6.19$ |
| Agent-R-8B | 63.91 | *N/A* | *N/A* | *N/A* |
| AgentLM-7B | 63.60 | $2.24 \pm 0.40$ | $0.28 \pm 0.07$ | $3.91 \pm 1.07$ |
| AgentLM-13B | 70.80 | $2.36 \pm 0.46$ | $0.30 \pm 0.08$ | $5.23 \pm 1.59$ |
| **MEM1-WebShop** | **70.87** | **$0.81 \pm 0.10$** | **$0.15 \pm 0.16$** | **$2.61 \pm 0.48$** |

Table 3: Performance comparison across environments for single-objective tasks. Arrows indicate the desired direction. (SFT) means training with SFT and applying MEM1's prompt and rollout. Note that DeepResearcher is specifically trained on the single-objective Online Web-QA task with F1 score as the optimization objective, and Search-R1 is specifically trained on the single-objective Wiki-RAG task with EM as the objective.

| Environment | System | EM ↑ | F1 ↑ | Peak Token ($\times 10^2$) ↓ | Dependency ($\times 10^5$) ↓ | Inference Time ↓ |
|---|---|---|---|---|---|---|
| Wiki RAG | Qwen2.5-7B-Inst (truncate) | 0.287 | 0.382 | $6.28 \pm 0.05$ | $1.65 \pm 0.04$ | $2.26 \pm 0.04$ |
| | Qwen2.5-7B-Inst (A-MEM) | 0.246 | 0.373 | $8.47 \pm 0.12$ | $0.92 \pm 0.03$ | $11.2 \pm 0.40$ |
| | Qwen2.5-7B-Inst | 0.269 | 0.390 | $9.32 \pm 0.19$ | $1.17 \pm 0.04$ | $2.31 \pm 0.04$ |
| | Qwen2.5-14B-Inst | 0.422 | **0.534** | $8.89 \pm 0.21$ | $2.22 \pm 0.10$ | $6.73 \pm 0.24$ |
| | Search-R1 | **0.445** | 0.516 | $11.0 \pm 0.25$ | $1.50 \pm 0.05$ | **$2.23 \pm 0.14$** |
| | DeepResearcher | 0.419 | 0.503 | $13.3 \pm 0.34$ | $7.04 \pm 0.33$ | $3.86 \pm 0.09$ |
| | **MEM1-QA (SFT)** | 0.302 | 0.358 | $6.54 \pm 0.05$ | $3.30 \pm 0.13$ | $4.84 \pm 0.21$ |
| | **MEM1-QA** | 0.405 | 0.471 | **$5.63 \pm 0.03$** | **$0.76 \pm 0.02$** | $3.79 \pm 0.07$ |
| Online Web-QA | Qwen2.5-7B-Inst | 0.334 | 0.451 | $8.37 \pm 0.18$ | $1.39 \pm 0.06$ | $2.20 \pm 0.04$ |
| | DeepResearcher | 0.372 | **0.492** | $10.27 \pm 0.19$ | $2.86 \pm 0.14$ | $2.87 \pm 0.06$ |
| | **MEM1-QA** | **0.397** | 0.485 | **$5.79 \pm 0.06$** | **$0.44 \pm 0.02$** | **$1.84 \pm 0.03$** |

## 4.4 Analysis on Emergent Agent Behaviors

Through analyzing MEM1's multi-turn interaction traces trained on 2-objective QA, we observe a range of emergent behaviors that are critical for handling long-horizon, multi-objective tasks, demonstrating capabilities well beyond simple retrieval. First, MEM1 learns to **manage multiple questions concurrently** by maintaining a structured internal state. As shown in Fig. 5(a), when faced with two multi-hop questions, the agent stores and updates memory for each question separately, guiding subsequent searches based on the identified information gaps. In (b), MEM1 exhibits the ability to shift focus when progress on one question stalls, recognizing difficulty and prioritizing the more tractable objective. Meanwhile, MEM1 learns to **interleave reasoning and memory** in its internal state `<IS>`, weaving important information into its decision-making process to support both information retention and action selection. In Fig. 5 (c), MEM1 explicitly extracts important information from previous search results and leverages it to formulate the next query that best addresses the current information gap. In addition, (d) shows that when new, relevant information is retrieved, MEM1 explicitly reasons about its significance and selectively updates its memory. We believe that learning these interleaved behaviors is key to achieving efficiency gains in memory without degrading performance. Beyond behaviors unique to our multi-objective setup and memory architecture, MEM1 also exhibits several **general-purpose search strategies**. In (e), the agent performs self-verification, correcting an earlier misconception and issuing a new query for confirmation. In (f), complex queries are decomposed into manageable subgoals before initiating the search. In (g), for questions requiring multi-turn information gathering, MEM1 extracts key information from search results and uses it to inform the next search. In (h), when overly specific queries fail, MEM1 re-scopes its query to improve retrieval. Notably, many of these behaviors, including verification, making a plan, and iterative search, are also reported in recent studies on deep research agents [24, 69].
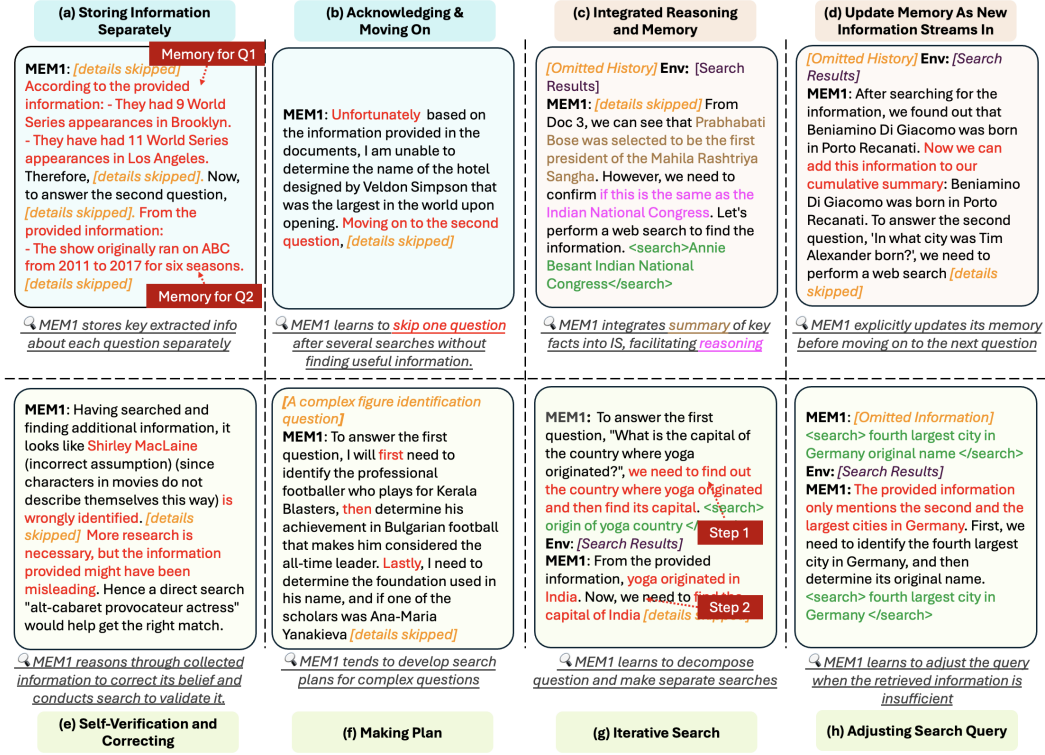
**(a) Storing Information Separately**

**MEM1:** *[details skipped]* According to the provided information: - They had 9 World Series appearances in Brooklyn. - They have had 11 World Series appearances in Los Angeles. Therefore, *[details skipped].* Now, to answer the second question, *[details skipped].* From the provided information: - The show originally ran on ABC from 2011 to 2017 for six seasons. *[details skipped]*

🔍 *MEM1 stores key extracted info about each question separately*

**(b) Acknowledging & Moving On**

**MEM1:** Unfortunately based on the information provided in the documents, I am unable to determine the name of the hotel designed by Veldon Simpson that was the largest in the world upon opening. Moving on to the second question, *[details skipped]*

🔍 *MEM1 learns to skip one question after several searches without finding useful information.*

**(c) Integrated Reasoning and Memory**

*[Omitted History]* **Env:** [Search Results]
**MEM1:** *[details skipped]* From Doc 3, we can see that Prabhabati Bose was selected to be the first president of the Mahila Rashtriya Sangha. However, we need to confirm if this is the same as the Indian National Congress. Let's perform a web search to find the information. <search>Annie Besant Indian National Congress</search>

🔍 *MEM1 integrates summary of key facts into IS, facilitating reasoning*

**(d) Update Memory As New Information Streams In**

*[Omitted History]* **Env:** *[Search Results]*
**MEM1:** After searching for the information, we found out that Beniamino Di Giacomo was born in Porto Recanati. Now we can add this information to our cumulative summary: Beniamino Di Giacomo was born in Porto Recanati. To answer the second question, 'In what city was Tim Alexander born?', we need to perform a web search *[details skipped]*

🔍 *MEM1 explicitly updates its memory before moving on to the next question*

**MEM1:** Having searched and finding additional information, it looks like Shirley MacLaine (incorrect assumption) (since characters in movies do not describe themselves this way) is wrongly identified. *[details skipped]* More research is necessary, but the information provided might have been misleading. Hence a direct search "alt-cabaret provocateur actress" would help get the right match.

🔍 *MEM1 reasons through collected information to correct its belief and conducts search to validate it.*

**(e) Self-Verification and Correcting**

*[A complex figure identification question]*
**MEM1:** To answer the first question, I will first need to identify the professional footballer who plays for Kerala Blasters, then determine his achievement in Bulgarian football that makes him considered the all-time leader. Lastly, I need to determine the foundation used in his name, and if one of the scholars was Ana-Maria Yanakieva *[details skipped]*

🔍 *MEM1 tends to develop search plans for complex questions*

**(f) Making Plan**

**MEM1:** To answer the first question, "What is the capital of the country where yoga originated?", we need to find out the country where yoga originated and then find its capital. <search> origin of yoga country </search> **Step 1**
**Env:** *[Search Results]*
**MEM1:** From the provided information, yoga originated in India. Now, we need to find the capital of India *[details skipped]* **Step 2**

🔍 *MEM1 learns to decompose question and make separate searches*

**(g) Iterative Search**

**MEM1:** *[Omitted Information]* <search> fourth largest city in Germany original name </search>
**Env:** *[Search Results]*
**MEM1:** The provided information only mentions the second and the largest cities in Germany. First, we need to identify the fourth largest city in Germany, and then determine its original name. <search> fourth largest city in Germany </search>

🔍 *MEM1 learns to adjust the query when the retrieved information is insufficient*

**(h) Adjusting Search Query**

Figure 5: Snippets of internal states and actions showing MEM1's Emergent Behaviors in 2-objective QA tasks. Light Blue denotes behaviors related to multi-objective tasks. Beige denotes behaviors related to memory in internal state. Pastel Green denotes behaviors related to general search strategies.

# 5 Conclusion, Limitations, and Future Work

We introduced MEM1, a reinforcement learning framework that enables language agents to perform long-horizon reasoning with consolidated memory. By integrating inference-time reasoning and memory consolidation into a unified internal state, MEM1 addresses the scalability challenges of prompt growth and achieves competitive performance across QA and web navigation benchmarks, with substantially reduced memory usage and inference latency. Despite these advantages, MEM1 assumes access to environments with well-defined and verifiable rewards. While this assumption holds in domains such as QA, math, and web navigation, many open-ended tasks present ambiguous or noisy reward structures. Fully realizing the potential of MEM1 therefore requires advances in modeling such tasks and designing suitable reward mechanisms—challenges that lie beyond the scope of this work. A promising future direction is to explore methods for training MEM1 agents in these open-ended settings where reward signals are sparse, delayed, or implicit.

# References

[1] A Cognitive Connection. Surprising exercises that will sharpen your short-term memory, January 2024. URL https://acognitiveconnection.com/surprising-exercises-that-will-sharpen-your-short-term-memory. Accessed: 2025-05-10.

[2] Chenxin An, Jun Zhang, Ming Zhong, Lei Li, Shansan Gong, Yao Luo, Jingjing Xu, and Lingpeng Kong. Why does the effective context length of llms fall short? In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.

[3] Anthropic. The claude 3 model family: Opus, sonnet, haiku. https://www.anthropic.com/news/claude-3-family, 2024.

[4] Alan D. Baddeley and Graham J. Hitch. Working memory. In Gordon H. Bower (ed.), *Psychology of learning and motivation*, volume 8, pp. 47–89. Academic Press, 1974.

[5] Hao Bai, Yifei Zhou, Jiayi Pan, Mert Cemri, Alane Suhr, Sergey Levine, and Aviral Kumar. Digirl: Training in-the-wild device-control agents with autonomous reinforcement learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 37:12461–12495, 2024.

[6] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.

[7] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.

[8] Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D. Lee, Deming Chen, and Tri Dao. Medusa: Simple LLM inference acceleration framework with multiple decoding heads. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 5209–5235, Vienna, Austria, 21–27 Jul 2024. PMLR. doi: 10.48550/arXiv.2401.10774. URL https://proceedings.mlr.press/v235/cai24b.html.

[9] Hyungjoo Chae, Namyoung Kim, Kai Tzu-iunn Ong, Minju Gwak, Gwanwoo Song, Jihoon Kim, Sunghwan Kim, Dongha Lee, and Jinyoung Yeo. Web agents with world models: Learning and leveraging environment dynamics in web navigation. *arXiv preprint arXiv:2410.13232*, 2024.

[10] Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. Accelerating large language model decoding with speculative sampling. *arXiv preprint arXiv:2302.01318*, February 2023. doi: 10.48550/arXiv.2302.01318. URL https://arxiv.org/abs/2302.01318.

[11] Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and Feng Zhao. Agent-flan: Designing data and methods of effective agent tuning for large language models. In *Findings of the Association for Computational Linguistics (ACL)*, pp. 9354–9366, 2024.

[12] Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Li YanTao, Jianbing Zhang, and Zhiyong Wu. Seeclick: Harnessing gui grounding for advanced visual gui agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 9313–9332, 2024.

[13] Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet Singh, and Deshraj Yadav. Mem0: Building production-ready ai agents with scalable long-term memory. *arXiv preprint arXiv:2504.19413*, 2025.

[14] ModelScope Community. SWIFT: A scalable lightweight infrastructure for fine-tuning. https://github.com/modelscope/ms-swift, 2024. Accessed: 2025-05-15.

[15] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, and et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

[16] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems (NeurIPS)*, 36:28091–28114, 2023.

[17] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. *arXiv preprint arXiv:2401.08281*, 2024.

[18] Google. Gemini: Try deep research and gemini 2.0 flash experimental. https://blog.google/products/gemini/google-gemini-deep-research/, 2024. Accessed: 2025-05-15.

[19] Xiangming Gu, Tianyu Pang, Chao Du, Qian Liu, Fengzhuo Zhang, Cunxiao Du, Ye Wang, and Min Lin. When attention sink emerges in language models: An empirical view. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.

[20] Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. A real-world webagent with planning, long context understanding, and program synthesis. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.

[21] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International Conference on Computational Linguistics (COLING)*, pp. 6609–6625, 2020.

[22] Jian Hu, Jason Klein Liu, and Wei Shen. Reinforce++: An efficient rlhf algorithm with robustness to both prompt and reward models. *arXiv preprint arXiv:2501.03262*, 2025. Version 3, revised 6 Apr 2025.

[23] Yucheng Jiang, Yijia Shao, Dekun Ma, Sina Semnani, and Monica Lam. Into the unknown unknowns: Engaged human learning through participation in language model agent conversations. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 9917–9955, 2024.

[24] Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*, 2025.

[25] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6769–6781, 2020.

[26] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019.

[27] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.

[28] Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative decoding. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 19274–19286, Honolulu, Hawaii, USA, 23–29 Jul 2023. PMLR. URL https://proceedings.mlr.press/v202/leviathan23a.html.

[29] Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. Compressing context to enhance inference efficiency of large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6921–6935, 2023.

[30] Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024.

[31] Zijun Liu, Peiyi Wang, Runxin Xu, Shirong Ma, Chong Ruan, Peng Li, Yang Liu, and Yu Wu. Inference-time scaling for generalist reward modeling. *arXiv preprint arXiv:2504.02495*, 2025.

[32] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems (NeurIPS)*, 36:46534–46594, 2023.

[33] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.

[34] Magnus Müller and Gregor Žunič. Browser use: Enable ai to control your browser. https://github.com/browser-use/browser-use, 2024.

[35] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2022.

[36] OpenAI. Gpt-4o system card. https://openai.com/index/gpt-4o-system-card/, 2024. Accessed: 2025-05-15.

[37] OpenAI. Introducing deep research, February 2025. URL https://openai.com/index/introducing-deep-research/.

[38] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. Measuring and narrowing the compositionality gap in language models. In *Findings of the Association for Computational Linguistics (EMNLP)*, pp. 5687–5711, 2023.

[39] Zehan Qi, Xiao Liu, Iat Long Iong, Hanyu Lai, Xueqiao Sun, Wenyi Zhao, Yu Yang, Xinyue Yang, Jiadai Sun, Shuntian Yao, et al. Webrl: Training llm web agents via self-evolving online curriculum reinforcement learning. *arXiv preprint arXiv:2411.02337*, 2024.

[40] Ruslan S. 4d masks support in transformers. https://huggingface.co/blog/poedator/4d-masks, 2024. Hugging Face Community Blog.

[41] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[42] Serper. Serper api: Fast and affordable google search api. https://serper.dev/, 2025. Accessed: 2025-05-15.

[43] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y.K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024.

[44] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv:2409.19256*, 2024.

[45] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 36:8634–8652, 2023.

[46] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.

[47] Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. Trial and error: Exploration-based trajectory optimization of llm agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 7584–7600, 2024.

[48] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, second edition, 2018. URL http://incompleteideas.net/book/the-book-2nd.html.

[49] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 12, pp. 1057–1063, 2000.

[50] OpenManus Team. Openmanus: Open-source ai agent framework. https://github.com/mannaandpoem/OpenManus, 2025.

[51] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, pp. 5998–6008, 2017.

[53] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(3):1–25, 2024.

[54] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*, 2022.

[55] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pp. 24824–24837, 2022.

[56] Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang, Kai-Wei Chang, and Dong Yu. Longmemeval: Benchmarking chat assistants on long-term interactive memory. In *International Conference on Learning Representations (ICLR)*, 2025.

[57] Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. A-mem: Agentic memory for llm agents. *arXiv preprint arXiv:2502.12110*, 2025.

[58] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

[59] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 2369–2380, 2018.

[60] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

[61] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. ReAct: Synergizing reasoning and acting in language models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023.

[62] Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. Lumos: Learning agents with unified data, modular design, and open-source llms. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2023.

[63] Chanwoong Yoon, Taewhoo Lee, Hyeon Hwang, Minbyul Jeong, and Jaewoo Kang. Compact: Compressing retrieved documents actively for question answering. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 21424–21439, 2024.

[64] Siyu Yuan, Zehui Chen, Zhiheng Xi, Junjie Ye, Zhengyin Du, and Jiecao Chen. Agent-r: Training language model agents to reflect via iterative self-training. *arXiv preprint arXiv:2501.11425*, 2025.

[65] Zhenrui Yue, Honglei Zhuang, Aijun Bai, Kai Hui, Rolf Jagerman, Hansi Zeng, Zhen Qin, Dong Wang, Xuanhui Wang, and Michael Bendersky. Inference scaling for long-context retrieval augmented generation. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.

[66] Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. Agenttuning: Enabling generalized agent abilities for llms. *arXiv preprint arXiv:2310.12823*, 2023.

[67] Jintian Zhang, Yuqi Zhu, Mengshu Sun, Yujie Luo, Shuofei Qiao, Lun Du, Da Zheng, Huajun Chen, and Ningyu Zhang. Lightthinker: Thinking step-by-step compression. *arXiv preprint arXiv:2502.15589*, 2025.

[68] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. Sglang: Efficient execution of structured language model programs. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.

[69] Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv preprint arXiv:2504.03160*, 2025.

[70] Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building autonomous agents. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.

[71] Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 43037–43067, 2023.

# A Details of MEM1

## A.1 Computing Resources and Training Details

All trainings of MEM1 are conducted on 4 H100 or H200 GPUs. We use the veRL framework [44] for RL and Swift [14] for SFT. For RL, both the data batch size and mini batch size are set to 64. Learning rate is set to $10^{-6}$ for the actor model and $10^{-5}$ for the critic model with a linear warmup of 50 steps. Temperature is set to 1 during training and 0.01 during inference.

All evaluations are conducted on a single H200 GPU, which serves the respective models as an API service using the vLLM framework [27] with automatic prefix caching enabled.

## A.2 RAG Configuration

For RAG on local Wiki corpus, we use Faiss-GPU [17] serving an E5 Base model [54]. The Wiki corpus is taken from a Wikipedia 2018 dump [25]. The number of passages for each retrieval is set to 3 for a fair comparison with other methods.

For online web search queries, we use Serper API [42], which offers Google search results including titles, snippets, and URLs. For each search, we return the top 10 results to the agent as external information. We do not ask the agent to retrieve the content of specific webpages.

## A.3 Prompts

> **Prompt 1: Multi-Objective Task (QA)**
>
> ```
> You will answer multiple complex questions using iterative reasoning,
> summarization, and web search.
>
> At each step, you will see the questions, a cumulative summary of
> relevant information, the current search query, and search results
> (except in the first step, where only the questions are provided).
> Your task is to:
>
> 1.  Perform reasoning and update a cumulative, concise summary
> within <think> ...  </think>.  This acts as persistent memory and
> must include all essential information from previous <think> and
> <information> tags.
>
> 2.  Then choose one of the following actions:
> - If any question remains unanswered, issue a single query for one
>   question inside <search> ...  </search>.  The query should consist of
>   keywords or a short phrase.  Only search one question at a time.
> - If all questions are answered, provide the final answers-separated
>   by semicolons-within <answer> answer1; answer2; ...  </answer>.  The
>   answers must be concise, contain only essential words, and avoid any
>   explanations.
>
> Important:
> - Always follow this structure after <information> or the initial
>   questions:  <think> ...  </think><search> ...  </search> or <think>
>   ...  </think><answer> ...  </answer>.
> - Do not search multiple queries or questions simultaneously.
>
> Answer the following questions:[QUESTIONS]
> ```

**Prompt 2: Single-Objective Task (QA)**

You will answer a complex question through iterative reasoning,
summarization, and web searches.

At each step, you can see the question, previous summary in <think>
...  </think>, search query in <search> ...  </search>, and the
returned information in <information> ...  </information> (except
the first step where you will be given only the question).  Then, you
should:

1.  Conduct reasoning, and then update a concise, cumulative summary
with essential information inside <think> </think>.  This is your
persistent memory and should include all important information from
previous and (i.e.
information and answers already found for questions).

2.  Then choose one:
- Issue a query (i.e., key words / phrases for search) inside(you may search repeatedly until the answer is clear).
This query will be used to conduct search and return the results in
<information> results </information>
- Provide the final concise answer (no explanations) if no additional
information is needed inside .  The answer should
be concise and only contain the words necessary to answer the
question.

After (or question at the beginning),
you should always follow the order:  <think> ...  </think><search> ...
</search> or <think> ...  </think><answer> ...  </answer>.

Question:  [QUESTION]

## A.4 Implementation Details of Metrics and Baselines

### A.4.1 Metrics

**Exact match.** In QA tasks, we use exact match (EM) as both the verifiable reward for the RL pipeline and the evaluation metric for the final output. The final response is extracted from between <answer> and </answer>. In multi-objective settings, the response should contain answers to each question separated by semicolons. If the XML tags are mismatched, or if the number of provided answers does not correspond to the number of questions, a score of $0$ is assigned. Otherwise, $1$ point is credited for each correct answer. During RL training, we do not provide any other intermediate rewards or format penalties, as we find that such manual interventions can interfere with the agent's learning process (see more in Sec. 4.4).

**F1 score.** The F1 score computes the harmonic mean between the precision $p$ and recall $r$. In the case of string matching, we split both the predicted answer and the ground truth. For example, if the ground truth is "United States of America", it is split into a list with lower-case words: "united", "states", "of", "america". The same works for the predicted answer. Then, denote the number of common words as $c$. Further denote the number of words in the predicted answer as $l$ and the number of words in the ground truth as $g$. Then, precision is calculated as $p := c/l$ and recall is calculated as $r := c/g$. The F1 score is finally computed as

$$\text{F1} := 2 \times \frac{p \times r}{p + r} .$$

If multiple ground truths are present, the maximum of all F1 scores is chosen. For multi-objective tasks, the final F1 is the sum of the F1 scores for each sub-question.

**Peak token usage.** Peak token usage is calculated as the maximum number of tokens (using GPT-4o-mini tokenizer) in any single sequence throughout the agent's entire trajectory. For fair comparison in our experiments, the system prompt is excluded when computing this sequence length. The peak token usage serves as a proxy for the inference-time memory requirement.

**Dependency length.** Following [67], the dependency metric is defined as the total number of historical tokens on which each generated token effectively depends. Let $T$ denote the total number of interaction steps. For each step $i \in [T]$, let $n_p^{(i)}$ be the number of prefix tokens and $n_o^{(i)}$ be the number of output tokens generated. The dependency metric is then calculated as

$$\texttt{Dependency} \coloneqq \sum_{i \in [T]} \frac{(2n_o^{(i)} + n_p^{(i)}) \times n_o^{(i)}}{2} \ .$$

At a high level, this metric quantifies the cumulative computational cost associated with the generation of an output trajectory. It is important to note that in MEM1, prefix tokens from previous steps are consolidated into a new internal state, rather than being continuously accumulated. In our experiments, we ignore the tokens in the system prompt when calculating the dependency metric.

**Inference time.** Inference time for each trajectory is recorded as the total elapsed time required to generate the complete output trajectory. For all experiments, these measurements are conducted on a single H200 GPU, operating with 10 concurrent threads. The vLLM inference framework is utilized, with its automatic prefix caching feature enabled.

### A.4.2 Baselines

**Search-R1.** As detailed in [24], the model is trained on the 1-objective task with the same dataset as MEM1. Search-R1 also uses exact match as its reward function. In comparison, MEM1 is trained exclusively on 2-objective tasks.

**Deep Researcher.** As detailed in [69], the model is trained on 1-objective task with a curated set from various QA datasets including HotPotQA and Natural Questions. Deep Researcher adopts the F1 score as the reward function.

### A.5 Algorithm

We provide an outline of the rollout of MEM1, which actively manages its context in Alg. 1. Parts of the pseudo-code follow [24].

**Algorithm 1** MEM1 Rollout

---

**Require:** Task prompt $x$, policy model $\pi_\theta$, world model $\mathcal{W}$, maximum turn $T$
**Ensure:** Final response $y$
 1: Initialize rollout sequence $y \leftarrow \varnothing$
 2: Initialize turn count $t \leftarrow 0$
 3: **while** $t < T$ **do**
 4:     Initialize current policy rollout sequence $y_t \leftarrow \varnothing$
 5:     **while** True **do**
 6:         Generate response token $y_r \sim \pi_\theta(\cdot \mid x, y + y_t)$
 7:         Append $y_r$ to rollout sequence $y_t \leftarrow y_t + y_r$
 8:         **if** $(t = T - 1)$ and $y_r \in [$`</answer>`, `<eos>`$]$ **then**
 9:           **break**                                 *// prevent the agent from searching further*
10:         **else if** $y_r \in [$`</query>`, `</answer>`, `<eos>`$]$ **then**
11:           **break**
12:         **end if**
13:     **end while**
14:     $y \leftarrow y_t$                                        *// all previous context removed.*
15:     **if** `<query>` `</query>` detected in $y_t$ **then**
16:         Extract search query $q \leftarrow \mathrm{Parse}(y_t,$ `<query>`, `</query>`$)$
17:         Retrieve environment feedback $d \leftarrow \mathcal{W}(q)$ from local storage, Search engine, HTML, $\cdots$
18:         `HINT` $\leftarrow$ `You have {`$T - t$`} turns left.`
19:         Insert $d$ into rollout $y \leftarrow y +$ `<info>`HINT$+ d$`</info>`
20:     **else if** `<answer>` `</answer>` detected in $y_t$ **then**
21:         **return** final generated response $y$
22:     **else**
23:         Mark the sample as invalid
24:     **end if**
25:     Increment turn count $t \leftarrow t + 1$
26: **end while**
27: **return** final generated response $y$

---

## A.6   MEM1 on Webshop Training Details

We use the same rollout pipeline and policy update mechanism for training MEM1 on WebShop. Compared to the QA tasks, we use a tailored prompt that retains the gist of memory consolidation with instructions specific to the WebShop environment, as shown in Prompt 3. Another distinction is that the WebShop environment comes with its own reward function corresponding to each state. Therefore, we do not use exact match but the built-in reward function as the reward signal when training in WebShop environment. The training and test splits also follow the original paper [60], with the first 1000 samples as the test set, the 1000th to 1500th as the val set, and the remaining as the train set.

## A.7   Additional Discussion on the Attention Matrix Design.

We wish to note that our modification to the attention matrix *does not* fully recover the attention of the original trajectories because of the change in position ids. Specifically, prior works [28, 10, 8] that utilized the attention matrix to compress multiple trajectories mainly targeted tree-exploration, i.e., generating multiple sequences with the same prefix. For these works, on top of the attention matrix, they adjusted the position ids as well, so each trajectory follows a consecutive increasing position ids. However, in MEM1, the prefix does not remain the same because of memory consolidation. This results in each `<IS>` having two possible position ids, one for the previous turn and one for the next turn. To completely recover the original attention, we need to duplicate each `<IS>` and assign different position ids to the two copies. However, such duplication can significantly slow down training because the training trajectories are now much longer.

As such, for training efficiency, we do not duplicate the `<IS>` and assign the position ids for the previous trajectory to each `<IS>`. While this modification slightly deviates from the "ideal"

implementation, effectively, it can be viewed as simply adding white spaces in the training trajectories and has no significant impact on the experimental results.

## B  Broader Impacts

MEM1 opens up the potential to enable more scalable, efficient, and intelligent AI agents capable of sustaining long, goal-directed interactions in dynamic environments. As AI systems are increasingly deployed in complex real-world tasks—such as scientific research, legal analysis, personalized education, and digital customer service—models must go beyond single-turn capabilities and manage evolving contexts over many steps. MEM1's memory-consolidation mechanism allows language models to maintain high performance without the growing computational and environmental costs typically associated with long-context processing. By reducing inference-time memory and compute demands, MEM1 paves the way for more sustainable and scalable AI deployment, making advanced reasoning agents accessible to a wider range of users and institutions, including those with limited resources. Moreover, MEM1's unified framework of reasoning and context consolidation sets a precedent for future research on intelligence that can learn to adapt, reflect, and summarize information autonomously, inspiring more trustworthy, interpretable, and human-aligned AI systems.
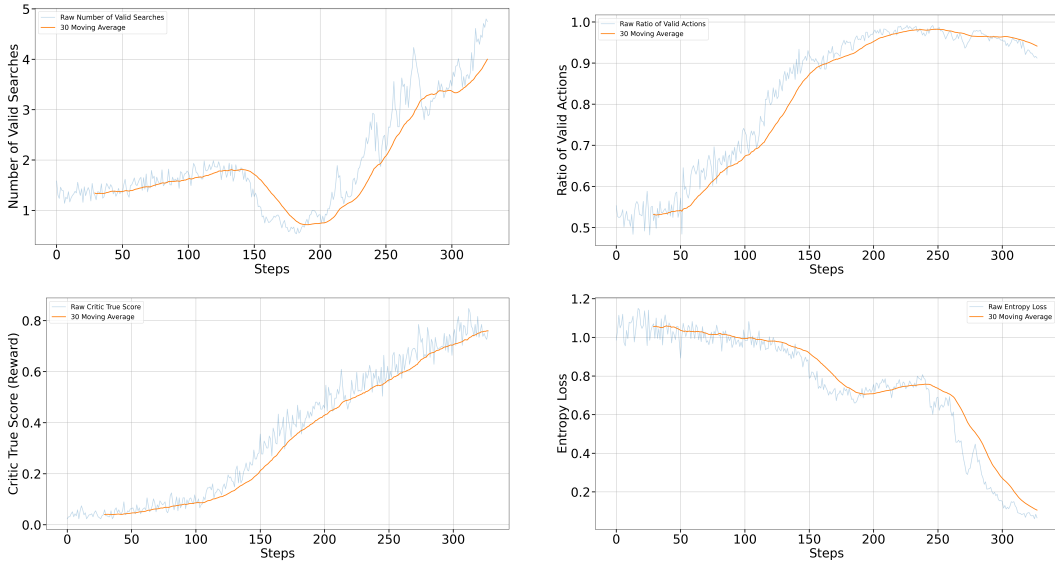
## C  Training Trajectory Analysis of MEM1



Figure 6: Metrics of training progresses for MEM1 with RL.

We present the training dynamics of the 2-objective QA-trained MEM1 in Fig. 6, where several distinct phases emerge during the learning process. In the initial exploration phase (first 50 steps), the agent demonstrates little task proficiency. The reward remains consistently low, while the entropy loss is high, suggesting random or undirected behavior. The ratio of valid actions hovers around 0.55, indicating that the agent frequently fails to follow the expected output format. During this period, MEM1 has not yet learned to reliably use the required structure involving `<query>` and `<answer>` tags.

Shortly after, we observe the onset of format acquisition. The agent gradually improves its structural consistency, reflected in the rising ratio of valid actions. This improved adherence to format correlates with an increase in reward, suggesting that proper formatting directly contributes to the agent's task success. By around step 150, a notable behavioral shift occurs. The number of valid searches begins to drop sharply, while the reward continues to increase. This implies that the agent has discovered a shortcut: by reducing the number of searches—perhaps to avoid format violations—it can maintain high format fidelity and improve its reward without fully solving the task. This short-horizon

21

optimization suggests the agent is exploiting the reward structure, favoring formatting compliance over content completeness.

Between steps 150 and 200, the agent enters a phase of refined format mastery. The ratio of valid actions steadily climbs, but the number of searches remains low. During this phase, reward growth slows, and entropy begins to flatten. The plateau in entropy indicates that the agent is looking for new policies to boost the reward. At this stage, the agent has reached a local optimum: it's producing valid but under-informed answers.

After step 200, a second behavioral shift occurs. The number of valid searches begins to rise again, suggesting that the agent is learning to extend its interaction horizon to gather more information. The agent learns to balance formatting constraints with information acquisition. As a result, the reward increases more sharply. Finally, after step 250, the agent enters a phase of policy consolidation. The entropy loss drops sharply—signaling a transition from exploration to exploitation—as the agent settles into a more deterministic, high-reward policy. By this stage, the agent effectively combines format compliance, sufficient searching, and high-quality answer generation.

# D    Analysis on Implementation Details

## D.1    RL Generalizes Better Than SFT

A natural question arises: can Supervised Fine-Tuning (SFT) with high-quality trajectories match the performance of reinforcement learning (RL)? To investigate this, we compare MEM1-QA trained via RL against MEM1-QA (SFT), where both models are trained on the 2-objective QA task. Additionally, the SFT model is further trained on 1-objective and 3-objective QA tasks to enhance its generalization ability. As shown in Tab. 4, the SFT model consistently underperforms compared to its RL counterpart across tasks with varying numbers of questions (objectives). Notably, when the number of objectives exceeds six, the performance of the SFT model collapses, whereas the RL-trained model continues to demonstrate strong robustness and scalability.

Table 4: Comparison of RL and SFT on increasing number of multi-turn questions. Exact match scores ↑ is better. Gap shows absolute difference. Red numbers show collapsed SFT behavior.

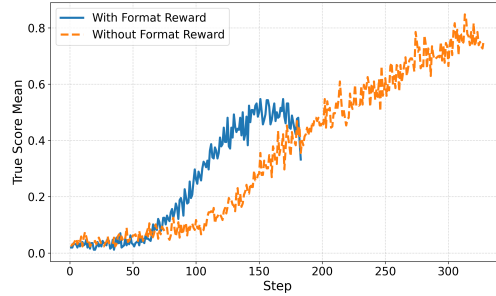| #Q | RL ↑ | SFT ↑ | Gap ↑ | RL Gain (%) ↑ |
|---|---|---|---|---|
| 1 | 0.410 | 0.300 | 0.110 | +36.7% |
| 2 | 0.709 | 0.433 | 0.276 | +63.7% |
| 3 | 0.976 | 0.648 | 0.328 | +50.6% |
| 4 | 1.120 | 0.626 | 0.494 | +78.9% |
| 6 | 1.630 | 0.088 | 1.542 | +1752% |
| 8 | 1.870 | 0.027 | 1.843 | +6826% |
| 16 | 1.900 | 0.000 | 1.900 | — |



Figure 7: Training curves comparing MEM1 trained with and without format reward.

## D.2    Format Reward Accelerates Convergence but Degrades Final Performance

It is common to incorporate format reward when training reasoning models and multi-turn reasoning agents [15, 69, 24]. In our study, we experimented with a format reward that enforces the agent to produce outputs using specific structural tags: `<IS>`, `<query>`, and `<answer>`. If the agent fails to use the expected tags correctly, the turn is terminated and a penalty of -1 is applied.

As shown in Fig. 7, using the format reward leads to faster convergence during training but results in worse final performance. The format-constrained agent achieves an exact match score of 0.466, compared to 0.709 for MEM1 trained with only outcome-based reward on the same testing set for the 2-objective QA task. Additionally, the format-constrained agent generates fewer tokens, with an average peak of 514.9 tokens, whereas the outcome-reward-trained MEM1 reaches an average peak of 640 tokens.

We hypothesize that the format reward accelerates structural learning but constrains exploration of effective reasoning strategies. As a result, the agent learns to produce shorter responses with valid syntax but develops less effective internal state representations, leading to degraded task performance.