# Multiple Memory Systems for Enhancing the Long-term Memory of Agent

**Gaoke Zhang[1], Bo Wang[1*], Yunlong Ma[1], Dongming Zhao[2], Zifei Yu[3]**

[1]College of Intelligence and Computing, Tianjin University
[2]AI Lab, China Mobile Communication Group Tianjin Co., Ltd
[3]Huizhi Xingyuan Information Technology Co., Ltd
{zhanggaoke, bo_wang}@tju.edu.cn

## Abstract

An agent powered by large language models have achieved impressive results, but effectively handling the vast amounts of historical data generated during interactions remains a challenge. The current approach is to design a memory module for the agent to process these data. However, existing methods, such as MemoryBank and A-MEM, have poor quality of stored memory content, which affects recall performance and response quality. In order to better construct high-quality long-term memory content, we have designed a multiple memory system (MMS) inspired by cognitive psychology theory. The system processes short-term memory to multiple long-term memory fragments, and constructs retrieval memory units and contextual memory units based on these fragments, with a one-to-one correspondence between the two. During the retrieval phase, MMS will match the most relevant retrieval memory units based on the user's query. Then, the corresponding contextual memory units is obtained as the context for the response stage to enhance knowledge, thereby effectively utilizing historical data. Experiments on LoCoMo dataset compared our method with three others, proving its effectiveness. Ablation studies confirmed the rationality of our memory units. We also analyzed the robustness regarding the number of selected memory segments and the storage overhead, demonstrating its practical value.

## Introduction

On the path towards Artificial General Intelligence (AGI), the development of large language models (LLMs) has achieved remarkable milestones, endowing them with robust language understanding and generation capabilities. Existing LLMs, such as ChatGPT (Brown et al. 2020), Gemini (Team et al. 2023), DeepSeek-r1 (Guo et al. 2025), Qwen (Yang et al. 2024), and Llama (Grattafiori et al. 2024), have driven breakthroughs in this field. LLM technologies are profoundly influencing the transformation of problem-solving paradigms across multiple AI domains.

Agents driven by LLM have emerged as an effective approach to addressing dialogue system challenges (Xi et al. 2025). The memory capabilities of agents are gradually becoming a pivotal factor in propelling them towards higher-level cognition and autonomous behavior (Zhang et al.

---

2024). Traditional short-term memory mechanisms are inadequate for meeting the demands of complex tasks, which require contextual understanding, long-term reasoning, and personalized responses (Wang et al. 2024). Consequently, constructing a memory system with durability, structure, and adaptability has become one of the core challenges in realizing truly "intelligent" agents.

To address this issue, researchers have explored both parametric and non-parametric approaches to optimize the performance of LLMs in long-form tasks. Parametrically, by augmenting LLMs with additional parameters, knowledge can be retained for future use (Wang et al. 2023). The Elastic Weight Consolidation algorithm (Huszár 2018) , developed through collaboration between DeepMind and Imperial College London, enables neural networks to retain knowledge from previous tasks while learning new ones, mitigating the "catastrophic forgetting" problem and marking a significant step towards continuous learning in AI. However, parametric memory is prone to generating distorted and non-factual outputs and lacks interpretability (Ji et al. 2023). Non-parametrically, external components are employed to enhance the long-term memory capabilities of LLMs. The external memory module of an agent can store historical dialogue content, often using the Retrieval-Augmented Generation (RAG) approach to store long-term memories as vectors in a database. When a new query arises, relevant vectors are matched from the vector database to serve as memory content (Gao et al. 2023). However, simply storing historical dialogues in a database does not adequately simulate the human memory formation process, nor does it effectively match information from historical dialogues.

Despite the proposal of various memory architectures, such as A-MEM (Xu et al. 2025), which extracts keywords, constructs summaries, and creates tags from dialogues as memory units, employing the Zettelkasten method to dynamically index and link knowledge networks, and MemoryBank (Zhong et al. 2024), which summarizes dialogue content and analyzes user personality and mood to serve as memory units, these approaches often fall short in practical scenarios. Users pose questions from diverse perspectives, yet the aforementioned methods simply extract keywords and summaries, saving them as memory content. This leads to a gap between users' queries and the content in memory units, with the memory content often lacking in qual-

ity, thereby affecting retrieval recall effectiveness. To enhance retrieval performance, we integrate Tulving's memory system classification (Tulving 1985) and the encoding specificity principle (Tulving and Thomson 1973), extracting episodic memory, semantic memory, cognitive perspectives, and keywords as our memory units content. This approach provides higher-quality memory content, thereby improving our retrieval recall effectiveness.

Currently, AI agents' memory systems encompass various types, including Short-Term Memory (STM), Long-Term Memory (LTM), Episodic Memory, Semantic Memory, and Procedural Memory (Sumers et al. 2023). STM is utilized for processing immediate contexts, such as the context window in dialogue systems. LTM, on the other hand, enables cross-session information storage and retrieval through databases, knowledge graphs, or vector embeddings (Gutiérrez et al. 2024; Xi et al. 2025; Hu et al. 2023). Episodic Memory allows agents to recall specific events, supporting case-based reasoning (Tulving et al. 1972). Semantic Memory stores structured factual knowledge, facilitating logical reasoning and knowledge retrieval (Tulving 1986). We have designed a multiple memory system based on cognitive psychology theory to generate high-quality long-term memory.

Our contributions are as follows:

(1) Current research rarely considers memory content quality, leading to its poor state. To address the issue of low-quality memory content in current memory systems, we are inspired by multiple memory systems theory, encoding specificity principle, and levels of processing theory to propose a multiple memory system to enhance the long-term memory capabilities of agents. The memory system extract keywords, multiple cognitive perspectives, episodic memory, and semantic memory as memory fragments through the analysis and processing of short-term memory to construct high-quality long-term memory content.

(2) In order to meet the requirements of different tasks in the retrieval stage and generation stage, we have built retrieval memory units for relevance matching with queries and contextual memory units as context for knowledge enhancement during generation.

(3) We conducted experiments on the LoCoMo dataset to evaluate the retrieval and generation capabilities, and also performed ablation study. These experimental results demonstrate the effectiveness of our method.

## Related Work

### Cognitive Psychology

The theory of multiple memory systems (Tulving 1985) posits that memory is not processed by a single system, but rather consists of multiple subsystems that are functionally independent and have different structural foundations. These systems are responsible for different types of information processing and storage, with specific neural bases and behavioral characteristics. Endel Tulving (Tulving et al. 1972) proposed a three-category classification model in 1985, including: Procedural Memory: involving the learning of skills and habits, such as riding a bicycle. Semantic Memory:

Stores factual and conceptual knowledge, such as the name of a capital city. Episodic memory: Recording personal experiences and events, such as the last birthday party. The Levels of Processing Theory (Craik and Lockhart 1972) suggests that the formation of memory does not depend on independent storage systems such as short-term memory and long-term memory, but rather on the depth and manner of information processing. The Encoding Specificity Principle (Tulving and Thomson 1973) states that the memory effect of information depends on its processing method and context during encoding. Specifically, when information is encoded, the environment, emotional state, sensory stimulation, and so on will all become part of the memory trace. Therefore, only when the conditions during retrieval match those during encoding, the retrieval of memory is most effective.

Inspired by the levels of processing theory, we process short-term memory content via a multi-level approach to form long-term memory representations. Considering the multiple memory theory, we categorize semantic and episodic memories as long-term memory components. Recognizing users' diverse questions, which reflect varied perspectives on the initial memory content they have, we are considering the encoding specificity principle to enhance recall by aligning coding content closely with the question's context. We integrate keywords from short-term memory and diverse cognitive angles derived from short-term into long-term memory representations to boost matching efficacy.

### Memory System

In recent years, research on the long-term memory mechanisms of agents has exhibited a diversifying trend. MemoryBank (Zhong et al. 2024) achieves long-term storage and retrieval of knowledge in multi-turn dialogues by introducing an external memory bank, integrating three modules—a writer, a retriever, and a reader—and leveraging the Ebbinghaus memory curve theory. It boasts strong scalability and modularity advantages, yet it suffers from issues such as coarse-grained memory content and suboptimal memory selection effects. MemoChat (Lu et al. 2023), on the other hand, injects user-related information into the model in the form of static summaries using concise, manually constructed "memos," significantly enhancing consistency and efficiency in open-domain dialogues. However, it relies on high-quality memo construction and lacks the capability for dynamic learning and memory adjustment. Think-in-Memory (Liu et al. 2023) proposes a "pre-recollection + post-reflection" framework that mimics human-like cognition, explicitly separating retrieval and integrated reasoning into distinct stages and introducing a self-reflection mechanism to enhance reasoning capabilities, making it suitable for complex tasks. Nevertheless, it incurs high reasoning costs, is sensitive to retrieval, and demands significant prompt engineering efforts. Integrating Dynamic Human-like Memory Recall and Consolidation in LLM-Based Agents (Hou, Tamoto, and Miyashita 2024) divides memory into short-term and long-term categories, dynamically updating and managing memories based on factors
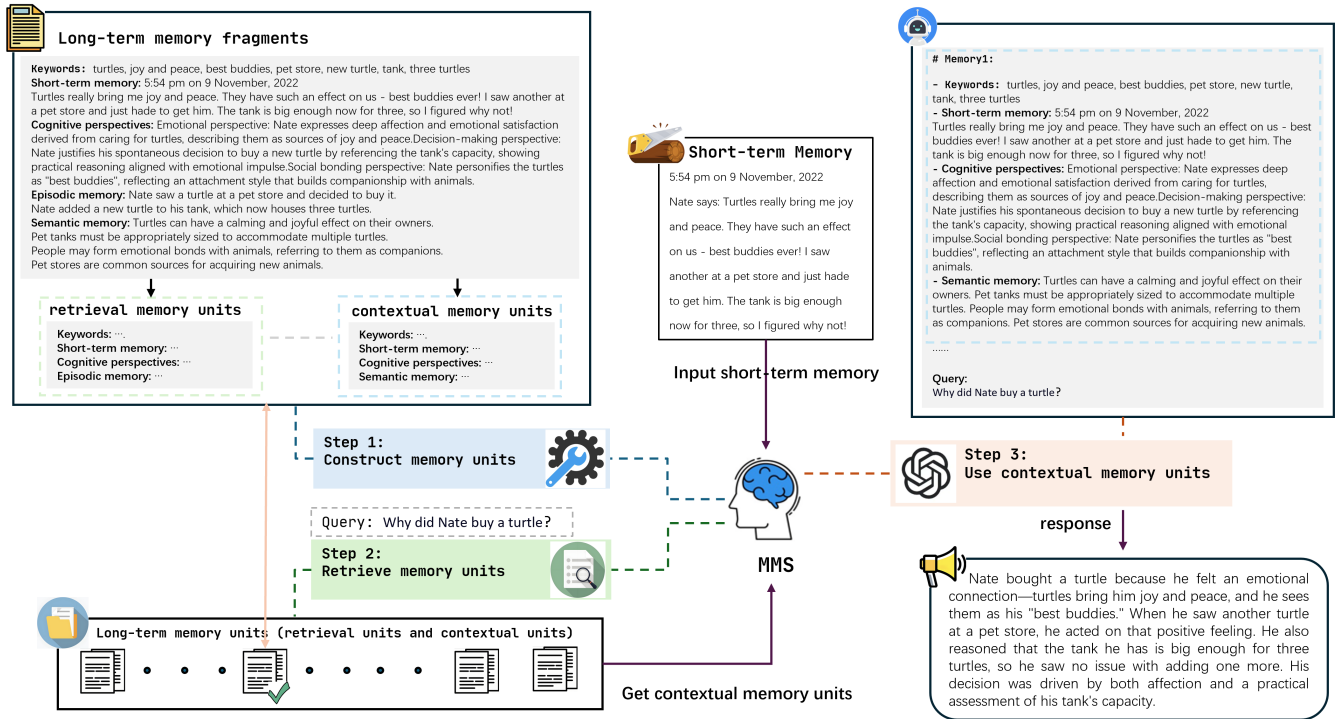
Figure 1: Schematic of multi-memory system process: After acquiring short-term memory, MMS processes it into memory fragments and constructs retrieval units and contextual memory units. During retrieval, the k most relevant retrieval units are matched to the query, and their corresponding contextual units are then used as context input for the agent's response.

such as usage frequency and temporal proximity. ChatDB (Hu et al. 2023) employs a database as the symbolic memory for LLMs. Its approach involves structuring textual data and storing it in a database, enabling LLMs to swiftly retrieve accurate knowledge through database queries when relevant information is needed. Human memory is characterized by features such as forgetting, consolidation, and association.

These methods lack detailed design for the most critical part of the memory module - the quality of the memory content, and only obtain basic information such as simple keywords and summaries as memory content. We believe that high-quality memory content is the key to improving recall and response abilities. We will combine cognitive psychology theories and focus on designing high-quality long-term memory segments to improve these two abilities.

## Multiple Memory Systems

### Overview

Short-Term Memory temporarily stores and processes limited info for current tasks, while Long-Term Memory stores info long-term with high capacity, housing our knowledge, experiences, and skills. The brain extracts key info from short-term memory, processing it into long-term storage. The levels of processing Theory suggests information processing affects memory. MMS will comprehensively process and analyze short-term memory to create high-quality long-term memory fragments.

The multi-memory fragment system in this paper mimics this process, using the content of a round of dialogue $C$ as short-term memory $M_{short}$, extracting key information $M_{key}$ from the dialogue content, and further processing and analyzing the original dialogue content to construct different cognitive perspectives $M_{cog}$, episodic memory $M_{epi}$, and semantic memory $M_{sem}$. Then, these memory fragments are used to construct retrieval memory units $MU_{ret}$ and contextual memory units $MU_{cont}$, which are used for memory retrieval and memory use, respectively.

Our multi-memory fragment system involves three processes: the construction of long-term memory units, the retrieval of long-term memory units, and the use of contextual memory units. The process of the system processing short-term memory into long-term memory is shown in Figure 1.

### The Construction of Long-term Memory Units

The process of converting short-term memory into long-term memory involves two parts. First, the information in short-term memory is processed to generate multiple types of memory fragments. Then, these memory fragments are used to construct retrieval and matching memory units for retrieval and matching, and contextual memory units for the context of LLM.

**Short-term Memory Processing** The memory system processes information into $M_{longterm}$, analyzes $M_{short}$ through LLM, and constructing $M_{key}$, $M_{cog}$, $M_{epi}$, and $M_{sem}$. First, keywords are extracted from $M_{short}$ as impor-

tant textual identification information in short-term memory. Cognitively analyze short-term memory from different perspectives, and construct a multi-dimensional cognition of the short-term memory to enhance the matching effect. Event information such as plot events is used as $M_{epi}$, and MMS are used to analyze factual information in short-term memory, constructing it as plot memory content. Taking knowledge points as factual knowledge, using cognitive modules to analyze the content of knowledge points in short-term memory, and constructing them as $M_{cog}$. The disclosure of fragmented memory and long-term memory is as follows:

$$M_{key}, M_{cog}, M_{epi}, M_{sem} = LLM(M_{short}) \quad (1)$$

$$M_{longterm} = (M_{key}, M_{short}, M_{cog}, M_{epi}, M_{sem}) \quad (2)$$

**Long-term Memory Fragment Storing** After completing information processing, we have five parts that constitute the long-term memory $M_{longterm}$: $M_{key}$, $M_{short}$, $M_{cog}$, $M_{epi}$ and $M_{sem}$. We will store long-term memory into two types, one for retrieval and matching, and the other for context-based knowledge enhancement. The keywords in long-term memory, the original short-term memory content, the various cognitive perspectives of short-term memory, and the sentence characteristics of episodic memory in short-term memory are more relevant to the semantic form of the user's query. However, the semantic memory of short-term memory is listed in the form of knowledge points, which is a higher-level knowledge extraction of short-term memory content. It may have some differences from the user's query language form and is suitable for user knowledge enhancement, but not for retrieval matching. Therefore, considering the principle of encoding specificity, We use $M_{key}$, $M_{short}$, $M_{cog}$ and $M_{epi}$ as retrieval memory units $MU_{ret}$, and construct them into vectors $V_{memory}$ for use in the retrieval matching stage.

Episodic memory describes event information, which is similar to the semantics of short-term memory information, and LLM have the ability to understand its content through short-term memory alone. Therefore, we believe that the content of episodic memory does not need to be input as contextual content into large models, only keywords, short-term memory, multiple cognitive perspectives, and semantic memory are required. Therefore, the composition of the retrieval memory units and the context memory units is as follows:

$$MU_{ret} = (M_{key}, M_{short}, M_{cog}, M_{epi}) \quad (3)$$

$$MU_{cont} = (M_{key}, M_{short}, M_{cog}, M_{sem}) \quad (4)$$

## Retrieval of Long-term Memory Units

The retrieval part will convert the user query $Q$ into a vector form, and then use the cosine similarity formula to calculate the similarity value between the user query vector $V_{query}$ and the constructed vector $V_{memory}$, selecting the top-k vectors $V_k = \{V_1, V_2, ..., V_K\}$ as the selected vectors. The cosine similarity used is disclosed as follows:

$$\cos\_sim(\mathbf{q}, \mathbf{v}) = \frac{\mathbf{q} \cdot \mathbf{v}}{\|\mathbf{q}\| \|\mathbf{v}\|} \quad (5)$$

where q and v are respectively the query vector and the vector stored in the memory system.

## Utilization of Long-term Memory Units

Based on the selection of top-k vectors $V_k = \{V_1, V_2, ..., V_K\}$, these retrieval memory units $MU_{ret}$ are mapped to corresponding contextual memory units $MU_{cont}$. $M_{key}$, $M_{short}$, $M_{cog}$, $M_{epi}$, and $M_{sem}$ are used as contextu $C_m$ inputs to LLM to response the user's query $Q$ and obtain a response $R$.

$$R = LLM(MU_{longterm}, Q) \quad (6)$$

# Experimental Setup

## Datasets

The LoCoMo dataset (Maharana et al. 2024) is designed to evaluate the long-term dialogue memory capability of large language models. Its primary task is to determine whether these models can accurately recall information from earlier dialogue after multiple rounds. The dataset contains 10 extended sessions, with each session averaging around 600 conversations and 26,000 tokens. Each conversation includes an average of 200 questions along with their corresponding correct answers. This dataset supports multiple evaluation scenarios and is currently a widely adopted authoritative dataset.

There are five problem types: (1) **Single-hop questions**: Test the memory system's basic retention by requiring accurate extraction of a single fact from long dialogue history. (2) **Multi-hop questions**: Assess the system's ability to integrate information across multiple dialogue rounds for reasoning. (3) **Temporal reasoning questions**: Evaluate if the system can understand event sequences and time evolution to construct a clear timeline. (4) **Open-domain knowledge questions**: Examine the system's recall and retrieval ability by combining dialogue context with external knowledge. (5) **Adversarial questions**: Challenge the system's resistance to forgetting and misleading by inserting interference information.

## Evaluation Metrics

We evaluate the performance of memory system using several standard metrics, including Recall@N (R@1, R@3, R@5), F1 score, and BLEU-1.

**Recall@N(R@N)** measures the average proportion of ground-truth answers retrieved within the top-N results. When the number of ground-truth items is less than N, the denominator is adjusted using $\min(N, |Gold_i|)$ to ensure fair evaluation. Formally:

$$\text{Recall}@N = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{|\text{Top-}N_i \cap \text{Gold}_i|}{\min(N, |\text{Gold}_i|)} \quad (7)$$

where $Q$ is the set of all evaluation queries, $i$ indexes a specific query, $\text{Gold}_i$ denotes the set of ground-truth answers

for the $i$-th query, $|\text{Gold}_i|$ is the number of ground-truth answers for that query, and Top-$N_i$ is the set of top-N results retrieved by the system for that query.

**F1 Score** is the harmonic mean of precision and recall. Let $TP$, $FP$, and $FN$ denote the number of true positives, false positives, and false negatives, respectively:

$$\text{Precision} = \frac{TP}{TP + FP} \tag{8}$$

$$\text{Recall} = \frac{TP}{TP + FN} \tag{9}$$

$$\text{F1} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \tag{10}$$

F1 is widely used in classification and span-level extraction tasks, such as named entity recognition or QA.

**BLEU-1** evaluates unigram precision between the generated output and the reference. BLEU-1 is computed as:

$$\text{BLEU-1} = \frac{\text{\# matching unigrams}}{\text{\# unigrams in hypothesis}} \tag{11}$$

A brevity penalty is typically applied to discourage overly short outputs. BLEU-1 is appropriate for tasks like machine translation and text generation where word-level overlap is informative.

## Baselines

We use three methods for comparison, asking: Naive RAG, MemoryBank, and A-MEM.

**Naive RAG** Only the character's dialogue content is vectorized as memory. For a question, it's vectorized and compared with all memory vectors for similarity. The top-k vectors of original dialogue are chosen as context for answering.

**MemoryBank** (Zhong et al. 2024) The memory system comprises three main methods: storage, which saves daily chats, event summaries, and user personality and mood assessments; retrieval, encoding dialogues and summaries into vectors for subsequent recall; and memory intensity update, utilizing an exponential decay model to mimic the Ebbinghaus forgetting curve.

**A-MEM** (Xu et al. 2025) The method has four steps: note construction (creating memory notes with conversation content, timestamps, keywords, tags, context, and links), link generation (identifying relevant historical memories, analyzing their ties to LLM for memory evolution), memory evolution (deciding on memory updates based on new notes), and memory retrieval (using cosine similarity between question and note vectors to fetch K most relevant memories for answers).

## Implementation Details

We employed GPT-4o, Qwen2.5-14B, and Gemini-2.5-pro-preview as base models, setting temperature to 0.5 for memory generation and 0.7 for question answering. For the latter, we input the top 5 most relevant contextual memory units (based on retrieval memory units) as context for the agent.

# Results and Discussion

We conducted an experimental comparative analysis on memory retrieval and memory use ability. In terms of memory retrieval, the recall rates (R@1, R@2, and R@5) were compared and analyzed. In terms of memory usage, F1, which measures the accuracy of responses, and BLEU-1, which measures the quality of responses, were compared and analyzed.

## Recall Results

We used GPT-4o, Qwen2.5-14B, and Gemini-2.5-pro-preview as the base models, and conducted experiments using our method and the other three methods for comparison, as shown in Table 1.

The experimental results show that MMS is superior to other methods in most cases, with a significant overall improvement effect. The improvement effect is most obvious in the Multi Hop scenario, achieving the maximum improvement across models and tasks in the three metrics of R@1, R@3, and R@5. GPT-4o and Gemini, in particular, improved recall by 8-11 points on the Multi Hop task compared to A-MEM. Multi Hop tasks rely on deep-level associative reasoning among multiple documents. The multi-memory fragment strategy of the MMS is more suitable for this complex retrieval scenario, enhancing the effectiveness of cross-document integration and reasoning. Significant improvements have also been observed in Open Domain and Temporal scenarios, indicating that MMS exhibits robustness in contexts with ambiguous and uncertain information structures. In the Single Hop scenario, most cases outperform other methods, but a few cases slightly underperform. This may introduce some redundancy in a simple MMS like Single Hop.

## Answer Results

MMS has achieved leading performance on multiple mainstream large language models, significantly outperforming existing methods such as NaiveRAG, MemoryBank, and A-MEM. Whether it is F1 or BLEU-1, MMS achieves comprehensive improvement in five types of tasks: single-hop, multi-hop, temporal, open-domain, and adversarial question answering. In particular, it has particularly outstanding advantages in multi-hop reasoning and open-domain tasks, demonstrating excellent information integration and reasoning depth. The strong performance of MMS may be due to its use of a multi-level memory structure, which has obvious advantages in improving question-answering accuracy, handling complex logical relationships, and resisting input interference. The experimental results are shown in Table 2.

## Ablation Study

We conducted ablation experiments on each module of the MMS. When the keyword part is not included, the R@1 metric under the Single Hop task and the R@5 metric under the Open Domain task perform the best. We believe that when the keyword part is not included. When multiple cognitive perspectives are not included, the R@3 metric under the Multi Hop task, the R@5 metric under the Temporal

| Model | Method | Single Hop | | | Multi Hop | | | Temporal | | | Open Domain | | | Adversarial | | | Average | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 |
| | NaiveRAG | 14.89 | 15.07 | 20.90 | 26.79 | 37.23 | 44.11 | 7.61 | 14.49 | 19.38 | 20.45 | 34.52 | 42.15 | 9.64 | 20.18 | 28.81 | 15.88 | 24.30 | 31.07 |
| GPT-4o | MemoryBank | 15.60 | 12.59 | 14.99 | 23.05 | 30.52 | 36.47 | 9.78 | 13.41 | 15.80 | 13.67 | 20.87 | 26.38 | 8.52 | 17.15 | 23.99 | 14.12 | 18.93 | 23.53 |
| | A-MEM | 24.82 | 23.76 | 29.73 | 33.02 | 49.79 | 58.96 | 16.30 | 22.28 | 30.02 | 29.01 | 44.81 | 53.90 | 10.54 | 20.96 | 25.11 | 22.74 | 32.32 | 39.54 |
| | **MMS** | **28.53** | **30.18** | **34.06** | **44.18** | **59.87** | **67.05** | **23.73** | **26.63** | **32.23** | **34.98** | **53.01** | **62.04** | **15.31** | **31.46** | **37.65** | **29.35** | **40.23** | **46.61** |
| | MemoryBank | 21.63 | 22.75 | 26.51 | 32.71 | 47.35 | 55.78 | 13.04 | 20.83 | 25.09 | 26.04 | 39.36 | 47.23 | 9.64 | 18.27 | 24.89 | 20.62 | 29.71 | 35.90 |
| Qwen2.5-14B | A-MEM | **25.21** | 25.71 | 27.45 | 39.19 | 52.49 | 59.01 | 21.73 | 23.19 | 32.46 | 31.27 | 46.63 | 54.82 | 14.57 | **27.58** | 32.85 | 26.39 | 35.12 | 41.32 |
| | **MMS** | 23.76 | **27.93** | 27.83 | **39.25** | **53.58** | **59.13** | **23.39** | **25.72** | **32.59** | 30.67 | **48.61** | **57.57** | **16.14** | 27.46 | **34.64** | 26.64 | 36.66 | 42.35 |
| | MemoryBank | 14.89 | 13.23 | 14.10 | 28.03 | 35.10 | 40.58 | 8.70 | 11.05 | 12.63 | 14.98 | 25.02 | 30.42 | 11.88 | 20.52 | 28.04 | 15.70 | 21.07 | 25.15 |
| Gemini-2.5-pro-preview | A-MEM | 17.73 | 17.43 | 22.42 | 25.55 | 38.32 | 46.73 | 10.87 | 15.04 | 19.93 | 22.47 | 37.06 | 43.32 | **13.46** | 23.99 | 29.93 | 18.02 | 26.37 | 32.47 |
| | **MMS** | **26.60** | **23.46** | **31.27** | **34.58** | **49.84** | **55.56** | **16.30** | **23.55** | **25.85** | **31.15** | **48.19** | **54.30** | 12.78 | **27.01** | **35.20** | **23.68** | **34.41** | **40.44** |

Table 1: The recall performance comparison of each method in five task scenarios. The gold truth corresponding to a single query may consist of multiple documents. For R@N, we use the minimum value between N and the number of documents in the gold truth as the denominator. Therefore, the denominators for R@1, R@3, and R@5 may differ, and it is reasonable to observe scenarios where, for example, R@3 has a lower value than R@1.

| Model | Method | Single Hop | | Multi Hop | | Temporal | | Open Domain | | Adversarial | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | F1 | BLEU-1 | F1 | BLEU-1 | F1 | BLEU-1 | F1 | BLEU-1 | F1 | BLEU-1 | F1 | BLEU-1 |
| GPT-4o | NaiveRAG | 18.04 | 12.24 | 33.03 | 28.28 | 14.03 | 12.15 | 31.74 | 26.78 | 7.60 | 6.90 | 20.89 | 17.27 |
| | MemoryBank | 15.39 | 11.18 | 28.44 | 24.18 | 11.91 | 11.23 | 22.04 | 18.86 | 8.65 | 7.93 | 17.29 | 14.68 |
| | A-MEM | 22.98 | 15.98 | 34.35 | 29.57 | 14.66 | 13.21 | 35.64 | 30.19 | 9.24 | 8.58 | 23.37 | 19.65 |
| | **MMS** | **28.67** | **20.96** | **47.37** | **39.98** | **20.81** | **19.07** | **42.98** | **36.89** | **12.87** | **11.67** | **30.54** | **25.74** |
| Qwen2.5-14B | NaiveRAG | 18.56 | 12.31 | 29.26 | 25.44 | 13.90 | 12.28 | 31.60 | 27.40 | 10.77 | 9.46 | 20.82 | 17.38 |
| | MemoryBank | 21.37 | 14.79 | 28.41 | 24.26 | 13.97 | 12.13 | 32.99 | 28.28 | 11.03 | 10.04 | 21.55 | 17.90 |
| | A-MEM | **30.71** | **23.44** | 32.46 | 28.55 | 14.08 | 13.97 | 43.02 | 37.69 | 11.38 | 10.11 | 26.33 | 22.75 |
| | **MMS** | 28.91 | 22.40 | **34.00** | **29.40** | **14.52** | **14.01** | **44.69** | **38.76** | **13.28** | **11.98** | **27.08** | **23.31** |
| Gemini-2.5-pro-preview | NaiveRAG | 23.80 | 17.27 | 34.96 | 26.92 | 14.24 | 12.50 | 30.64 | 26.43 | 5.76 | 5.15 | 21.88 | 17.65 |
| | MemoryBank | 14.12 | 9.30 | 34.15 | 29.36 | 8.55 | 6.48 | 25.82 | 22.17 | 4.62 | 3.92 | 17.45 | 14.25 |
| | A-MEM | 22.69 | 15.83 | 40.97 | 35.63 | 12.42 | 11.11 | 36.23 | 31.42 | **7.20** | **7.02** | 23.90 | 20.20 |
| | **MMS** | **25.96** | **18.96** | **41.49** | **36.99** | 12.91 | 11.32 | **42.00** | **36.59** | 7.09 | 6.81 | **25.89** | **22.13** |

Table 2: Comparison of Generation Performance of Four Methods on the LoCoMo Dataset

| Method | Single Hop | | | Multi Hop | | | Temporal | | | Open Domain | | | Adversarial | | | Avg | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 | R@1 | R@3 | R@5 |
| w/o Key | **30.85** | 29.72 | 34.02 | 42.68 | 59.29 | 66.17 | 23.21 | 23.55 | 31.86 | 34.60 | 52.55 | **62.18** | 11.88 | 27.36 | 33.40 | 28.64 | 38.49 | 45.53 |
| w/o Cog | 26.24 | 26.06 | 32.73 | 43.61 | **59.96** | 66.33 | 20.65 | **27.53** | 31.83 | 31.98 | 49.36 | 59.53 | **15.47** | 30.49 | 37.57 | 27.59 | 38.68 | 45.60 |
| w/o Epi | 27.66 | 27.30 | 33.17 | 42.37 | 59.03 | 68.38 | 18.47 | 24.99 | 31.68 | 31.51 | 47.98 | 57.51 | 15.02 | 29.48 | 37.00 | 27.01 | 37.76 | 45.55 |
| w/o Cog & EPi | 21.99 | 20.57 | 24.21 | 37.69 | 55.40 | 62.12 | 13.04 | 21.20 | 23.82 | 26.75 | 41.52 | 48.79 | 13.00 | 26.01 | 35.99 | 22.49 | 32.94 | 38.99 |
| MMS | 28.53 | **30.18** | **34.06** | **44.18** | 59.87 | **67.05** | **23.73** | 26.63 | **32.23** | **34.98** | **53.01** | 62.04 | 15.31 | **31.46** | **37.65** | **29.35** | **40.23** | **46.61** |

Table 3: The ablation experiment using GPT-4o as the base model of MMS method. The symbol "w/o" indicates an experiment in which a specific module was removed. Key represents the key words of short-term memory, Cog represents multiple cognitive perspectives on short-term memory, and Epi represents episodic memory of short-term memory.

task, and the R@1 metric under the Adversarial task perform the best. We believe that different memory fragments can be used to construct retrieval units based on different character scenarios. In a broader general scenario, the retrieval units composed of keywords, short-term memory, multiple cognitive perspectives, and episodic memory can achieve the best performance in most tasks. The experimental results are shown in Table 3.

## Analyze Long-term Memory Units

We used GPT-4o as the base to compare the effects of adding the remaining memory fragments in the MMS retrieval units and contextual memory units. For retrieval, retrieval memory units includes keywords, short-term memory, multiple cognitive perspectives, and episodic memory. We add semantic memory for comparison, and find recall rate decreases. This is because generated semantic memory is a high-level overview of factual knowledge in short-term memory, creating a gap with the semantic form of user's

| Method | Single Hop | | Multi Hop | | Temporal | | Open Domain | | Adversarial | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | F1 | BLEU-1 | F1 | BLEU-1 | F1 | BLEU-1 | F1 | BLEU-1 | F1 | BLEU-1 | F1 | BLEU-1 |
| w/o Key | 27.51 | 18.53 | 44.24 | 36.52 | 19.28 | 16.12 | 41.27 | 35.53 | 10.88 | 10,12 | 28.64 | 23.36 |
| w/o Cog | 27.77 | 19.77 | 45.68 | 38.46 | 17.36 | 15.47 | **43.01** | 36.87 | 12.27 | 11.65 | 29.22 | 24.44 |
| w/o Sem | 27.79 | 19.16 | 45.31 | 37.94 | 19.51 | 18.50 | 42.69 | 36.60 | 12.45 | 11.44 | 29.55 | 24.73 |
| w/o Cog&Sem | 28.24 | 20.24 | 45.05 | 38.34 | 17.92 | 15.44 | 42.71 | **37.42** | 12.83 | 11.53 | 29.35 | 24.59 |
| **MMS** | **28.67** | **20.96** | **47.37** | **39.98** | **20.81** | **19.07** | 42.98 | 36.89 | **12.87** | **11.67** | **30.54** | **25.74** |

Table 4: The ablation experiment using GPT-4o as the base model of MMS method. The symbol "w/o" indicates an experiment in which a specific module was removed. Key represents the key words of short-term memory, Cog represents multiple cognitive perspectives on short-term memory, and Sem represents semantic memory of short-term memory.
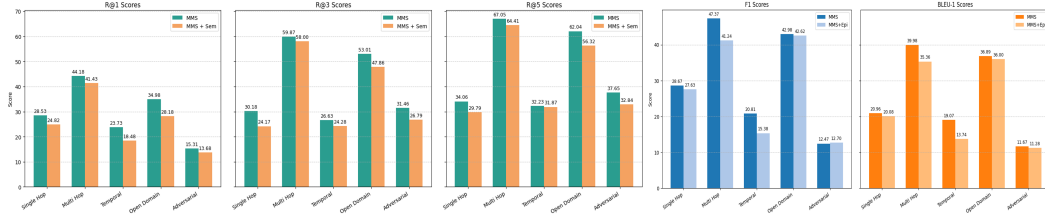


Figure 2: Compare the impact on performance after adding other segments. In terms of recall metrics, MMS and MMS+Sem were compared. In terms of generation, MMS and MMS+Epi were compared. Sem refers to semantic memory, and Epi refers to episodic memory.

question. For generation, contextual memory units of MMS consists of keywords, short-term memory, multiple cognitive perspectives, and semantic memory. We add contextual memory for comparison and find generation quality decreases after adding episodic memory. As episodic memory describes event-related content already present in original short-term memory, adding more can lead to redundancy and lower generation quality. Thus, the composition modules of retrieval and contextual memory units we use are reasonable. The experimental results are shown in Figure 2.

### Robustness Analysis of the Number of Memories

Table 5: Performance comparison for different values of n

| Metrics | n=1 | n=3 | n=5 | n=7 | n=9 |
|---|---|---|---|---|---|
| Avg Fl | 20.74 | 25.07 | 30.54 | 34.81 | 36.13 |
| Avg BLUE-1 | 17.44 | 21.39 | 25.74 | 28.95 | 31.28 |

We manipulated the quantities of memory fragments to evaluate the impact on performance. As the value of n increased, performance exhibited enhancement notwithstanding the introduction of noise, owing to the high-caliber memory contents facilitating the differentiation between pivotal and noisy data. This demonstrates the resilience of our method with respect to the quantities of memory fragments.

### Analysis of Token and Latency Overhead

We analyzed token and latency overhead for different methods experimentally, specifically calculating the average overhead needed to generate the memory content for

Table 6: Comparison of Latency Overhead and Token Overhead

| Metrics | MMS | A-MEM | Memory Bank |
|---|---|---|---|
| Avg Latency | 1.309 | 3.931 | 0.949 |
| Avg Tokens | 744 | 1429 | 238 |

each query. Compared to A-MEM, our approach proves to be faster and more resource-efficient. Although there's a slight rise in latency and overhead compared to the memory repository, its long-term memory is overly simplistic and of low quality, resulting in poor performance. Conversely, our method generates a larger volume of high-quality memory content despite the increased overhead, with only a minimal latency increase that barely affects user experience. Our method holds practical value.

### Conclusion

This paper creates a multi-memory system by integrating with cognitive psychology to build effective long-term memory, boosting recall and generation quality. Multiple memory theory suggests human memory comes in many forms. Current methods don't consider the variety of human memory fragments and just use summarization, leading to low-quality content. Since people understand problems from different angles, we view short-term memory's various cognitive aspects as long-term memory fragments. Levels-of-processing theory states that deeper encoding leads to better memory retention. MMS turns short-term memory into long-term fragments like keywords, different cogni-

tive views, episodic and semantic memories. It constructs retrieval units from these for recall and builds contextual units to enhance knowledge. Experiments indicate MMS enhances recall and generation at lower cost, proving practical value. Our analysis under varying memory segment counts reveals high-quality content ensures sustained gains and noise resilience. Due to space constraints, our approach centers on memory content, with future work to include more memory operation designs. Our approach effectively integrates cognitive psychology into the research on agent memory within artificial intelligence for future studies.

# References

Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.

Craik, F. I.; and Lockhart, R. S. 1972. Levels of processing: A framework for memory research. *Journal of verbal learning and verbal behavior*, 11(6): 671–684.

Gao, Y.; Xiong, Y.; Gao, X.; Jia, K.; Pan, J.; Bi, Y.; Dai, Y.; Sun, J.; Wang, H.; and Wang, H. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2: 1.

Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Vaughan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Guo, D.; Yang, D.; Zhang, H.; Song, J.; Zhang, R.; Xu, R.; Zhu, Q.; Ma, S.; Wang, P.; Bi, X.; et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Gutiérrez, B. J.; Shu, Y.; Gu, Y.; Yasunaga, M.; and Su, Y. 2024. Hipporag: Neurobiologically inspired long-term memory for large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

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

Hu, C.; Fu, J.; Du, C.; Luo, S.; Zhao, J.; and Zhao, H. 2023. Chatdb: Augmenting llms with databases as their symbolic memory. *arXiv preprint arXiv:2306.03901*.

Huszár, F. 2018. Note on the quadratic penalties in elastic weight consolidation. *Proceedings of the National Academy of Sciences*, 115(11): E2496–E2497.

Ji, Z.; Lee, N.; Frieske, R.; Yu, T.; Su, D.; Xu, Y.; Ishii, E.; Bang, Y. J.; Madotto, A.; and Fung, P. 2023. Survey of hallucination in natural language generation. *ACM computing surveys*, 55(12): 1–38.

Liu, L.; Yang, X.; Shen, Y.; Hu, B.; Zhang, Z.; Gu, J.; and Zhang, G. 2023. Think-in-memory: Recalling and post-thinking enable llms with long-term memory. *arXiv preprint arXiv:2311.08719*.

Lu, J.; An, S.; Lin, M.; Pergola, G.; He, Y.; Yin, D.; Sun, X.; and Wu, Y. 2023. Memochat: Tuning llms to use memos for consistent long-range open-domain conversation. *arXiv preprint arXiv:2308.08239*.

Maharana, A.; Lee, D.-H.; Tulyakov, S.; Bansal, M.; Barbieri, F.; and Fang, Y. 2024. Evaluating very long-term conversational memory of llm agents. *arXiv preprint arXiv:2402.17753*.

Sumers, T.; Yao, S.; Narasimhan, K.; and Griffiths, T. 2023. Cognitive architectures for language agents. *Transactions on Machine Learning Research*.

Team, G.; Anil, R.; Borgeaud, S.; Alayrac, J.-B.; Yu, J.; Soricut, R.; Schalkwyk, J.; Dai, A. M.; Hauth, A.; Millican, K.; et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Tulving, E. 1985. How many memory systems are there? *American psychologist*, 40(4): 385.

Tulving, E. 1986. Episodic and semantic memory: Where should we go from here? *Behavioral and Brain Sciences*, 9(3): 573–577.

Tulving, E.; and Thomson, D. M. 1973. Encoding specificity and retrieval processes in episodic memory. *Psychological review*, 80(5): 352.

Tulving, E.; et al. 1972. Episodic and semantic memory. *Organization of memory*, 1(381-403): 1.

Wang, L.; Ma, C.; Feng, X.; Zhang, Z.; Yang, H.; Zhang, J.; Chen, Z.; Tang, J.; Chen, X.; Lin, Y.; et al. 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6): 186345.

Wang, W.; Dong, L.; Cheng, H.; Liu, X.; Yan, X.; Gao, J.; and Wei, F. 2023. Augmenting language models with long-term memory. *Advances in Neural Information Processing Systems*, 36: 74530–74543.

Xi, Z.; Chen, W.; Guo, X.; He, W.; Ding, Y.; Hong, B.; Zhang, M.; Wang, J.; Jin, S.; Zhou, E.; et al. 2025. The rise and potential of large language model based agents: A survey. *Science China Information Sciences*, 68(2): 121101.

Xu, W.; Liang, Z.; Mei, K.; Gao, H.; Tan, J.; and Zhang, Y. 2025. A-mem: Agentic memory for llm agents. *arXiv preprint arXiv:2502.12110*.

Yang, A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Li, C.; Liu, D.; Huang, F.; Wei, H.; et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.

Zhang, Z.; Bo, X.; Ma, C.; Li, R.; Chen, X.; Dai, Q.; Zhu, J.; Dong, Z.; and Wen, J.-R. 2024. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501*.

Zhong, W.; Guo, L.; Gao, Q.; Ye, H.; and Wang, Y. 2024. Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 19724–19731.