

# BMAM: Brain-inspired Multi-Agent Memory Framework

Yang Li<sup>1</sup> Jiaxiang Liu<sup>1</sup> Yusong Wang<sup>2</sup> Yujie Wu<sup>3</sup> Mingkun Xu<sup>1\*</sup>

<sup>1</sup>Guangdong Institute of Intelligence Science and Technology, Zhuhai, China

<sup>2</sup>Institute of Science Tokyo <sup>3</sup>The Hong Kong Polytechnic University

{liyong, liujiaxiang, xumingkun}@gdiist.cn

wangyi@lr.pi.titech.ac.jp yu-jie.wu@polyu.edu.hk

## Abstract

Language-model-based agents operating over extended interaction horizons face persistent challenges in preserving temporally grounded information and maintaining behavioral consistency across sessions, a failure mode we term *soul erosion*. We present **BMAM** (Brain-inspired Multi-Agent Memory), a general-purpose memory architecture that models agent memory as a set of functionally specialized subsystems rather than a single unstructured store. Inspired by cognitive memory systems, BMAM decomposes memory into episodic, semantic, salience-aware, and control-oriented components that operate at complementary time scales, organised as a six-phase memory lifecycle. To support long-horizon reasoning, BMAM organises episodic memories along explicit timelines and retrieves evidence by fusing multiple complementary signals. Experiments on the **LoCoMo** benchmark show that BMAM achieves **78.45% accuracy**, outperforming seven memory-augmented baselines. Pairwise ablations reveal *super-additive* synergy between brain-region components rather than redundant stacking, and a Soul Portability Test demonstrates **87.5% identity-integrity** across full memory export, clear, and restore. A targeted refinement of the temporal-trigger heuristics raises Long-MemEval multi-session accuracy from 45.2% to 56.4%, validating the architectural decomposition behind BMAM. Code is available at <https://github.com/innovation64/BMAM>.

## 1 Introduction

Language-model-based agents increasingly operate in settings that require maintaining and reasoning over information accumulated across extended interactions, spanning diverse tasks, domains, and time scales. Such agents must retain past experiences, organize them into usable memory struc-

tures, and retrieve relevant information under varying goals and contexts. However, large language models are constrained by finite context windows and lack an explicit mechanism for managing long-term memory beyond the current input (Packer et al., 2023; Maharana et al., 2024). Retrieval-augmented generation (RAG) partially alleviates this limitation by fetching external documents on demand, but it treats memory as an external text repository rather than an internal, evolving system. As a result, RAG-style approaches provide limited support for persistent memory accumulation, temporal organization, and cross-episode reasoning, motivating the need for a general-purpose memory framework that can support long-horizon agent behavior across tasks rather than task-specific retrieval pipelines (Zhang et al., 2025b; Hu et al., 2025).

Evidence from cognitive science suggests that memory is not a single monolithic store, but is supported by multiple functionally specialized subsystems operating over complementary time scales (e.g., fast episodic encoding alongside slower semantic consolidation and executive control) (O’Reilly et al., 2014). Inspired by this view, we propose **BMAM** (Brain-inspired Multi-Agent Memory Framework), a brain-inspired multi-agent memory architecture that decomposes agent memory into interacting subsystems responsible for episodic storage, semantic consolidation, salience-aware selection, and intent-conditioned control (Li et al., 2025b). BMAM constructs internal memory representations rather than relying solely on external retrieval, and employs a timeline-indexed episodic memory organization to support temporally grounded access to past experiences. The framework further integrates a hybrid retrieval mechanism that combines lexical, dense, knowledge-graph, and temporal signals via reciprocal rank fusion, together with asynchronous memory consolidation processes inspired by comple-

\*Corresponding author.

mentary learning principles. To coordinate memory access across different temporal scales, BMAM adopts a hierarchical memory control mechanism from recent work, enabling both fast context-level access and slower consolidated memory retrieval. In preliminary analyses of long-horizon agent behavior, we observe a recurring failure pattern in which fragmented or misaligned memory leads to degradation in temporal coherence and identity-related behavior across interactions, which we refer to as soul erosion, providing a diagnostic lens for failures of long-term memory management in general-purpose agent settings. Our main contributions are:

- We identify and characterize **soul erosion**, a recurring failure pattern in long-horizon agent behavior where fragmented or misaligned memory leads to degradation in temporal coherence and identity-related behavior.
- We propose **BMAM**, a brain-inspired framework that addresses this challenge by decomposing memory into specialized subsystems (episodic, semantic, salience). Crucially, we introduce a **timeline-indexed organization** and a **hybrid retrieval strategy** that fuses lexical, semantic, and temporal signals for robust grounding.
- We validate BMAM on the **LoCoMo** benchmark, achieving **78.45%** accuracy and outperforming baselines in long-horizon settings. Further ablation studies empirically confirm the critical role of the hippocampus-inspired subsystem in enabling temporal reasoning.

**Soul Erosion: Why Memory Matters** We use the term *soul erosion* to describe a recurring failure pattern in long-horizon agent interactions, where fragmented or misaligned memory leads to degradation in behavioral continuity and identity-related behavior. Analogous to how human identity relies on the continuity of autobiographical memory (Wilson and Ross, 2003; Bluck and Liao, 2013), an AI agent’s “soul” (its consistent preferences, behavioral tendencies, and interaction patterns) may gradually degrade when long-term memory is poorly organized or inconsistently accessed.

**Formal Definition** We formalize soul erosion as a composite degradation metric over three orthogonal dimensions. Let  $\mathcal{M}_t$  denote the agent’s

memory state at interaction step  $t$ . We define the **soulfulness score**  $\mathcal{S}$  as:

$$\mathcal{S}(\mathcal{M}_t) = \alpha \cdot T(\mathcal{M}_t) + \beta \cdot C(\mathcal{M}_t) + \gamma \cdot I(\mathcal{M}_t) \quad (1)$$

where  $T(\cdot)$  measures *temporal coherence* (ability to correctly order and recall when events occurred),  $C(\cdot)$  measures *semantic consistency* (absence of factual contradictions), and  $I(\cdot)$  measures *identity preservation* (retention of user-specific preferences and traits). The weights  $\alpha, \beta, \gamma \geq 0$  with  $\alpha + \beta + \gamma = 1$  reflect task-specific importance.

**Soul erosion** is then defined as the degradation of soulfulness over time:

$$\mathcal{E}(t_0, t) = \mathcal{S}(\mathcal{M}_{t_0}) - \mathcal{S}(\mathcal{M}_t) \quad (2)$$

where  $t_0$  is a reference point (e.g., initial interaction or last memory consolidation). A positive  $\mathcal{E}$  indicates soul erosion has occurred. In our experiments, we operationalize these components using benchmark proxies:  $T$  via LoCoMo temporal accuracy,  $C$  via cross-session consistency metrics, and  $I$  via PrefEval and PersonaMem scores.

**Soul erosion encompasses three distinct failure modes** (Figure 1), each arising from different memory failures and requiring specialized countermeasures:

**(1) Temporal Erosion** The agent loses track of *when* events occurred, leading to anachronistic or temporally inconsistent responses. Cognitive research shows that temporal context is fundamental to episodic memory organization (Howard and Kahana, 2002; Eichenbaum, 2014), and benchmarks like LoCoMo and LongMemEval (Maharana et al., 2024; Wu et al., 2025) reveal that LLM agents frequently fail on temporal queries. As shown in Figure 1 (left), without explicit temporal organization, the agent may confuse event order, overlook durations, or fail to answer time-dependent queries. BMAM addresses temporal erosion through StoryArc timeline indexing, which maintains explicit temporal structure over stored experiences.

**(2) Semantic Erosion** Facts and relationships degrade or become internally inconsistent across interactions. This mirrors the forgetting and interference phenomena studied in human memory (Wixted, 2004; Anderson, 2003), where memories compete and degrade without proper consolidation. As depicted in Figure 1 (center), the agent may provide contradictory answers about the same entity over time. HippoRAG (Gutiérrez et al.,

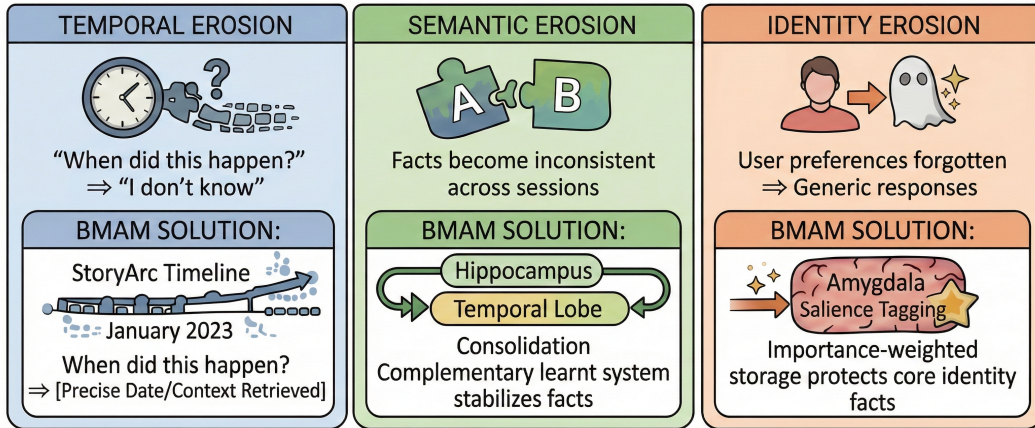


Figure 1: Soul erosion types and BMAM countermeasures. Each erosion mechanism requires a specialized defense: temporal erosion is addressed by StoryArc timeline indexing, semantic erosion by hippocampus-to-temporal-lobe consolidation, and identity erosion by amygdala salience tagging.

2024) and memory surveys (Zhang et al., 2025b) highlight this challenge. BMAM counters semantic erosion through hippocampus-to-temporal-lobe consolidation, which promotes frequently accessed and high-confidence episodic memories into stable semantic representations.

**(3) Identity Erosion** User preferences, personality traits, and persistent behavioral patterns may be overwritten or lost as new context accumulates. Research on autobiographical memory emphasizes that identity coherence depends on preserving self-relevant experiences (Conway, 2005; McAdams, 2001). Benchmarks like PersonaMem and PrefEval (Jiang et al., 2025; Zhao et al., 2025) demonstrate that current systems struggle to maintain user-specific information. As shown in Figure 1 (right), this failure mode undermines personalization: the agent “forgets” who the user is. BMAM mitigates identity erosion through amygdala-inspired salience tagging, which prioritizes identity-relevant information and protects it from being overwhelmed by transient context.

### Multi-Agent Coordination as Erosion Defense

A central design insight of BMAM is that these three forms of erosion arise from distinct memory failures and cannot be fully addressed by a single mechanism. Cognitive neuroscience research demonstrates that human memory relies on multiple specialized systems (the hippocampus for episodic encoding, the neocortex for semantic consolidation, and the amygdala for emotional salience) that interact to maintain coherent long-term memory (O’Reilly et al., 2014). Inspired by

this functional specialization, BMAM distributes memory functions across multiple interacting components, each targeting a specific erosion type (Figure 1). Our ablation studies (Table 9) empirically validate this design: removing the hippocampus-inspired episodic memory causes the largest performance drop, confirming its critical role, while other components contribute complementary defenses against different erosion types.

## 2 Background and Related Work

**Memory Architectures for LLM Agents** RAG improves factual grounding but treats memory as implicit and transient; retrieved passages are not reorganized into stable internal structures (Zhang et al., 2025a; Wang et al., 2024b; Xu et al., 2024). Agent-centric frameworks address this limitation through explicit memory management. MemGPT pioneered virtual context management by treating LLMs as operating systems with hierarchical memory tiers (Packer et al., 2023). MemoryBank extends this with forgetting mechanisms inspired by Ebbinghaus curves (Zhong et al., 2024). A-MEM introduces agentic memory that autonomously manages storage and retrieval (Xu et al., 2025b). Production systems like Mem0, Memobase, and MemOS provide scalable memory APIs with multi-component stores (Chhikara et al., 2025; memodbio, 2025; Li et al., 2025a). Hierarchical approaches organize memory by semantic abstraction levels (Sun et al., 2026; Wang et al., 2025; Xu et al., 2025a). Memory-augmented transformers have explored various mechanisms for extending context, including segment-level recurrence (Dai et al.,

2019), kNN-augmented attention (Wu et al., 2022), and brain-inspired episodic memory (Das et al., 2024). These approaches primarily target token-level or sequence-level prediction, whereas BMAM targets long-horizon agent memory management: what to store, how to organize it temporally, and how to retrieve it under changing goals.

### Brain-Inspired and Cognitive Approaches

Cognitive neuroscience motivates separating fast episodic encoding from slower semantic consolidation and salience-based prioritization (O’Reilly et al., 2014). This principle has inspired systems like HippoRAG for hippocampus-style indexing (Gutiérrez et al., 2024), Nemori for event segmentation (Nan et al., 2025), and reflective memory systems that learn from experience through prospective and retrospective reflection (Tan et al., 2025; Shinn et al., 2023). Recent architectures emphasize tight coupling between perception and memory, forming closed loops that support adaptive long-term memory (Wang et al., 2024a; Park et al., 2023). Compared to these approaches, BMAM differs in three key aspects: (1) *multi-region coordination*: while HippoRAG focuses on hippocampal pattern separation, BMAM models interactions among multiple brain-region analogs (hippocampus, temporal lobe, amygdala, prefrontal cortex); (2) *explicit temporal indexing*: unlike Nemori’s event boundaries, BMAM maintains continuous timeline structures that support arbitrary temporal queries; (3) *salience-aware consolidation*: BMAM integrates amygdala-inspired importance signals into the consolidation process, prioritizing identity-relevant information over transient context.

**Benchmarks** Long-term memory benchmarks evaluate temporal reasoning (LoCoMo (Maharana et al., 2024), LongMemEval (Wu et al., 2025)), preference consistency (PrefEval (Zhao et al., 2025)), and persona recall (PersonaMem (Jiang et al., 2025)), providing complementary perspectives on the challenges BMAM addresses.

### 3 BMAM Framework

BMAM adopts a coordinator-centered multi-agent architecture that decomposes long-term memory into functionally specialized components while maintaining a unified memory substrate. A central coordinator routes information among interacting subsystems responsible for memory storage,

### Algorithm 1 BMAM Memory Lifecycle Pipeline

---

**Require:** query  $q$ , context  $C$ , regions  $\mathcal{B} = \{\text{Hip}, \text{Tmp}, \text{Amg}, \text{Pfc}, \text{Bg}\}$   
**Ensure:** response  $r$

```

# Phase 1: Query analysis & routing
1:  $f \leftarrow \text{ANALYZEFEATURES}(q, C)$  ▷ temporal/entity/intent
2:  $\langle b_1, \dots, b_K \rangle, s \leftarrow \text{LEARNABLEROUTER}(q)$ 
3:  $A \leftarrow \text{THALAMUSGATING}(q, C)$  ▷  $A[b] \in \{0, 1\}$ 
# Phase 2: Region-specific parallel access
4: for all activated region  $b \in \mathcal{B}$  in parallel do
5:    $M_{\text{hip}} \leftarrow \text{Hip.EPISODICSEARCH}(q)$ 
6:    $M_{\text{tmp}} \leftarrow \text{Tmp.SEMANTICSEARCH}(q)$  ▷ BM25 + vec + KG
7:    $M_{\text{amg}} \leftarrow \text{Amg.SALIENCERETRIEVE}(q)$ 
8:    $M_{\text{pfc}} \leftarrow \text{Pfc.WMLOOKUP}(q)$ 
9:    $M_{\text{bg}} \leftarrow \text{Bg.PROCEDURALMATCH}(q)$ 
10: end for
# Phase 3: Cross-region resonance scoring
11:  $M \leftarrow M_{\text{hip}} \cup M_{\text{tmp}} \cup M_{\text{amg}} \cup M_{\text{pfc}} \cup M_{\text{bg}}$ 
12: for all  $m \in M$  do
13:    $m.\rho \leftarrow |\{b: m \in M_b\}|$  ▷ resonance count
14:    $m.\sigma \leftarrow \sum_i w_i \cdot \phi_i(m, q)$  ▷ weighted score
15: end for
# Phase 4: Hippocampal–Prefrontal loop
16:  $i \leftarrow 0$ 
17: repeat
18:    $M \leftarrow \text{Amg.MODULATESALIENCE}(M, q)$ 
19:    $M \leftarrow M \cup \text{Tmp.FILLSEMANTICGAPS}(q, M)$ 
20:    $\gamma, \text{gaps} \leftarrow \text{Pfc.EVALUATEEVIDENCE}(q, M)$ 
21:   if  $\gamma < \tau$  then
22:      $q' \leftarrow \text{Pfc.GAPANALYSIS}(q, \text{gaps})$ 
23:      $M \leftarrow M \cup \text{Hip.EXPANDEDRETRIEVAL}(q')$ 
24:   end if
25:    $i \leftarrow i + 1$ 
26: until  $\gamma \geq \tau$  or  $i = I_{\text{max}}$ 
# Phase 5: Fusion, reranking, structured reasoning
27:  $M^* \leftarrow \text{LLMRERANK}(\text{TOPN}(M, \sigma), q)$ 
28: if  $\text{KGTRIGGER}(q)$  then  $M^* \leftarrow M^* \cup \text{Tmp.KGQUERY}(q)$ 
29: end if
30: if  $\text{TEMPQUERY}(q)$  then  $r_t \leftarrow \text{TEMPORALREASON}(q, M^*)$ 
31: end if
# Phase 6: Generation & memory write-back
32:  $r \leftarrow \text{GENERATE}(q, M^*, C)$ 
33:  $\text{Hip.ENCODEEPISODE}(q, r, C)$ 
34:  $\text{Amg.TAGEMOTION}(q, r)$ 
35: if  $\text{importance} > 0.6 \wedge \text{access} > 2$  then
36:    $\text{Tmp.CONSolidATE}(\text{Hip.episodes})$ 
37: end if
38:  $\text{Pfc.UPDATEWM}(q, r)$ 
39:  $\text{FEEDBACKLOOP.UPDATEVECTORS}(M^*, \gamma)$ 
40: return  $r$ 

```

---

retrieval, consolidation, and control, enabling modular specialization without fragmenting memory state.

**Memory Loop and Coordination** BMAM implements an explicit memory loop inspired by hippocampus–neocortex dynamics. Incoming experiences are encoded into episodic memory using fast, discriminative representations and tagged with salience signals, while relevant content is maintained in a constrained working-memory buffer to support immediate reasoning. Over time, selected episodic information is consolidated into semantic memory and a shared knowledge graph. Retrieval closes the loop by jointly accessing episodic and semantic evidence under temporal constraints, with feedback signals adjusting consolidation priorities and routing decisions.

### Functionally Specialized Memory Components

BMAM decomposes memory into complemen-

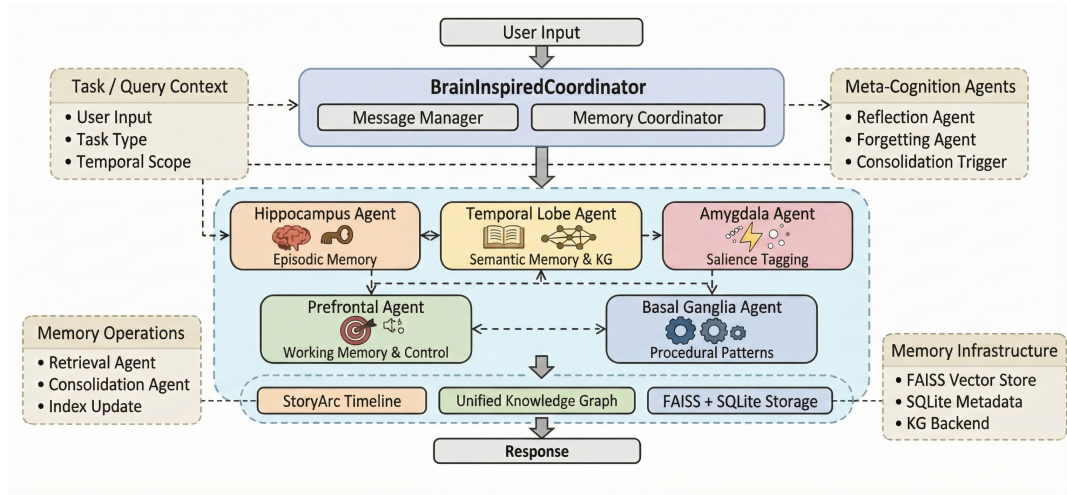


Figure 2: BMAM architecture overview. A central coordinator orchestrates multiple functionally specialized memory subsystems sharing a unified memory substrate with episodic timelines, a knowledge graph, and vector-based storage.

tary subsystems with explicit roles and capacities. Episodic memory stores temporally grounded interaction traces and supports discriminative addressing. Semantic memory consolidates stable facts and relations into a shared knowledge graph. A salience-aware component computes importance signals from interaction cues (e.g., novelty, conflict, or user feedback) that modulate consolidation scheduling and retrieval weighting.

The **Prefrontal** component implements executive control functions inspired by the prefrontal cortex’s role in working memory maintenance and cognitive control (Miller and Cohen, 2001). Specifically, it performs three functions: (1) *query routing*, classifying incoming queries along dimensions (temporal, identity, preference, factual) to determine which memory subsystems to consult; (2) *working-memory buffering*, maintaining a capacity-limited buffer (10 items) of recent context for immediate reasoning without full memory retrieval; and (3) *attention allocation*, dynamically weighting evidence sources based on query requirements. Control-oriented components, including the Prefrontal buffer and Basal Ganglia procedural patterns, together provide complementary protections against different forms of memory degradation.

**Unified Memory Substrate and Temporal Indexing** BMAM employs a unified memory substrate that combines key–value episodic storage, vector-based similarity indexing, and a shared knowledge graph. Episodic memories are organized into a timeline-indexed structure that records minimal narrative units indexed by entities, events, and times-

tamps. This temporal organization enables queries involving order, duration, and temporal relations (e.g., before/after, first/last), while consolidation processes selectively lift episodic information into semantic form to ensure consistency across representations.

**Hierarchical Coordination and Retrieval** To support long-horizon interactions, BMAM adopts hierarchical memory coordination mechanisms that regulate memory access and updates across multiple time scales. Fast paths support immediate context-level access, while slower paths govern semantic consolidation and procedural stabilization. Retrieval integrates fast-path detection, iterative interaction between episodic and control components, and uncertainty-driven multi-round retrieval. Evidence from episodic memory, semantic memory, and the knowledge graph is combined with temporal constraints, and feedback signals dynamically reweight lexical, dense, entity-based, and temporal cues.

**Memory Lifecycle** BMAM models memory as a dynamic lifecycle governing encoding, consolidation, retrieval, and revision, organised as a six-phase pipeline (Algorithm 1) and summarised diagrammatically in Appendix Figure 4.

**Input Analysis and Episodic Encoding** We first analyze each incoming interaction to extract entities, temporal expressions, and intent cues relevant to memory formation. The interaction is encoded as an episodic memory trace, capturing the contextual content together with inferred temporal and

semantic attributes. Saliency signals are computed from interaction cues (e.g., novelty, conflict, or user feedback) and attached to the episode. To support efficient short-term reasoning, a compact summary of recent episodes is maintained in a constrained working-memory buffer.

### Consolidation and Temporal Organization

Next, BMAM employs a complementary learning process in which frequently accessed and high-confidence episodic memories are selectively consolidated into semantic memory. Consolidated information populates a shared knowledge graph that maintains stable facts and relations across interactions. In parallel, episodic memories are organized into a timeline-indexed structure that records entity-centric events with associated temporal information. This temporal organization enables reasoning over event order, relative timing, and durations, supporting queries such as *when*, *before/after*, and *how long* without requiring full episodic recall.

### Hybrid Retrieval and Temporally Grounded Answering

To answer a query, BMAM retrieves relevant evidence from multiple sources, including episodic memory, semantic memory, and the timeline-indexed event structure. Each source  $s \in \mathcal{S}$  produces a ranked list of candidates, and lexical, dense, relational, and temporal signals are fused using (weighted) reciprocal rank fusion:

$$\text{score}(d | q) = \sum_{s \in \mathcal{S}} \frac{w_s}{k + \text{rank}_s(d | q)}, \quad (3)$$

where  $\text{rank}_s(d | q)$  is the rank of candidate  $d$  under source  $s$ ,  $k = 60$  is the smoothing constant following standard RRF practice, and  $w_s$  reflects the current preference over evidence sources. For time-dependent questions, temporal evidence is extracted from the timeline organization to compute relative orderings and durations, which are then used to generate temporally grounded answers. This retrieval process is adaptive: uncertainty and saliency signals may trigger additional retrieval rounds or reweight evidence sources.

### Background Optimization and Memory Revision

In parallel with online interaction, memory organization in BMAM is continuously refined through background processes. Episodic memories may be reconsolidated when re-accessed, increasing their stability or updating their content as new evidence emerges. Low-value or outdated memories are gradually pruned, while saliency-relevant

episodes receive prioritized consolidation. These processes allow BMAM to revise memory over time, preventing uncontrolled growth and reducing the accumulation of inconsistent or obsolete information.

### Computational Complexity

Local processing (routing, BM25/FAISS retrieval, resonance scoring) runs in tens of milliseconds; the hippocampal-prefrontal loop is bounded by  $I_{\max} = 3$ . End-to-end latency is dominated by LLM calls: **3–5 per query** (2–3 on simple queries; 4–5 when Pfc triggers re-retrieval). Measured on LoCoMo across eight runs, BMAM reaches **17.3 s/query** best and **22.4 s/query** average, at  $\sim \$0.001/\text{query}$  using gpt-4o-mini—comparable to rerank-RAG (2–3 calls) on simple queries, with additional cost concentrated on multi-hop queries, where BMAM’s margin over baselines is largest (Appendix Table 5).

### Continual Learning and Plasticity

Over longer interaction horizons, BMAM treats memory as a plastic substrate rather than a static store. Continual learning emerges from ongoing consolidation and reconsolidation, whereby retrieved evidence can update semantic memory instead of being frozen after first storage. Conceptually, if  $p_t(f)$  denotes the confidence of a semantic fact  $f$  at time  $t$ , and  $\hat{p}_t(f)$  is an evidence-based estimate from new retrieval/verification, then memory revision can be expressed as an exponential moving average:

$$p_{t+1}(f) = (1 - \lambda)p_t(f) + \lambda\hat{p}_t(f), \quad (4)$$

where  $\lambda \in (0, 1)$  is the update rate. This enables knowledge updates while damping noisy evidence. When confidence is low or information is incomplete, the system may actively seek clarification through follow-up interaction, strengthening memory traces and reducing uncertainty. Over time, adaptive routing, saliency-weighted storage, and confidence-calibrated retrieval (e.g., by adjusting  $w_s$  in Eq. 3) change what is stored, how it is indexed, and how evidence is combined, enabling BMAM to evolve its memory behavior as experience accumulates.

## 4 Experiments

### 4.1 Experimental Setup

We evaluate BMAM on four benchmarks designed to test long-horizon memory and personalization

capabilities (Table 1). Our evaluation focuses primarily on **LoCoMo** and **LongMemEval**, which together capture complementary challenges in conversational memory, temporal reasoning, and memory consistency over extended interactions.

**Primary Benchmarks** We focus primarily on **LoCoMo** (Maharana et al., 2024) and **LongMemEval** (Wu et al., 2025), which together capture complementary challenges in long-horizon memory. LoCoMo evaluates recall of facts, relationships, and events across extended multi-session dialogues, emphasizing single-hop factual recall, multi-hop reasoning, and temporally grounded questions. LongMemEval complements this with cross-session recall, preference tracking, knowledge updates, and explicit temporal reasoning, stressing memory consistency under evolving information.

**Additional Benchmarks** We further evaluate BMAM on **PersonaMem** and **PrefEval**, which focus on persona consistency and preference alignment, respectively. These benchmarks test whether memory systems can preserve user-specific information and behavioral preferences across interactions, complementing the conversational and temporal challenges posed by LoCoMo and LongMemEval.

Dataset	Scale	Task Focus	Metric
LoCoMo	10 groups, 1986 QA	long-horizon dialogue	Accuracy
LongMemEval	500 questions	long-term memory	Accuracy
PersonaMem	20 users, 589 QA	persona recall (MCQ)	Accuracy
PrefEval	1000 questions	preference alignment	Pers. rate

Table 1: Datasets and evaluation metrics used in BMAM experiments.

**Baselines** We compare BMAM against seven memory-augmented LLM systems: **MemOS** (Li et al., 2025a), a memory operating system with unified memory scheduling; **Mem0** (Chhikara et al., 2025), a scalable memory-centric architecture with optional graph-based memory; **MIRIX** (Wang and Chen, 2025), a multi-agent system with six specialized memory types; **Zep** (Rasmussen et al., 2025), a temporally-aware knowledge graph engine; **Memobase** (memodb-io, 2025), **Supermemory** (supermemoryai, 2025), and **MemU** (NevaMind-AI, 2025). Baseline results are

from Li et al. (2025a); we re-run MemOS with GPT-4o-mini for fair comparison.

**Evaluation Protocol** For all benchmarks, persistent memory is reset between independent evaluation units (e.g., LoCoMo conversation groups or individual users) while being preserved within each unit to reflect realistic interaction histories. Conversation logs are ingested through BMAM’s memory lifecycle prior to evaluation, and queries are issued in evaluation mode without additional learning. We follow the evaluation protocol, metrics, and judge prompts from MemOS<sup>1</sup>; baselines are evaluated using their official scripts. Crucially, to ensure a fair comparison, we re-evaluated the strongest baseline (MemOS) using the identical LLM backend (GPT-4o-mini) as BMAM, eliminating discrepancies arising from model version updates. During evaluation, all background processes (consolidation, reconsolidation, pruning) are disabled; memory state is frozen after ingestion to prevent test-time learning.

## 4.2 Results

Across benchmarks, BMAM is strongest on LoCoMo and PrefEval, while PersonaMem remains challenging: its multiple-choice format requires exact surface-form matching, whereas BMAM’s retrieval is optimized for open-ended generation, and its emphasis on shallow persona attributes differs from BMAM’s focus on temporally grounded identity. We provide a more detailed discussion in Appendix A.5. Temporal reasoning remains a key open challenge, and improving normalized temporal outputs and cross-session integration is an important direction for future work. Table 3 summarizes our main results.

**Temporal-Trigger Refinement** An audit of the temporal reasoning module showed that broad trigger heuristics (bare keywords *before*, *after*, *when*) and a low confidence threshold (0.35) caused over-firing on non-temporal queries. Three targeted fixes—narrowing triggers to precise forms (*when did*, *how long ago*), raising the threshold to 0.55, and expanding the non-temporal filter set from 4 to 30+ patterns—raise temporal accuracy from 76.9% to 84.6% and yield 78.89% overall on a single LoCoMo group ( $N=199$ ), consistent with the 10-group main result (78.45%); full per-category numbers are in Appendix Table 11.

<sup>1</sup>Official repository: <https://github.com/MemTensor/MemOS>; paper: <https://arxiv.org/abs/2507.03724>

Method	Single-hop	Multi-hop	Temporal	Open-domain	Overall
<b>BMAM (ours)</b>	<b>82.00</b>	<b>70.42</b>	62.31	79.55	<b>78.45</b>
MemOS-1031 <sup>†</sup>	64.54	57.29	71.34	79.90	73.90
Memobase	73.12	64.65	<b>81.20</b>	53.12	72.01
Mem0	73.33	58.75	52.54	45.83	64.57
MIRIX	68.32	54.26	68.54	46.88	64.33
Zep	65.23	52.12	54.82	33.33	59.22

Table 2: Head-to-head LoCoMo accuracy per category against seven memory-augmented baselines. BMAM leads the re-run MemOS baseline by 4.55 points overall, 8.67 on single-hop, and 5.77 on multi-hop. Memobase is stronger on temporal-only queries; this is addressed by the trigger refinement (Table 11). <sup>†</sup>Re-run with gpt-4o-mini for strict comparability; other baselines use reported numbers.

Dataset	Metric	Score	Correct/Total
LoCoMo	Accuracy	<b>78.45%</b>	1558/1986
LongMemEval	Accuracy	67.60%	338/500
PersonaMem	Accuracy	48.9%	288/589
PrefEval	Pers. rate	72.90%	729/1000

Table 3: BMAM results across four long-term memory benchmarks.

**LoCoMo Performance** On LoCoMo, BMAM achieves **78.45%** overall under our MemOS-aligned evaluation protocol. Table 2 reports the full head-to-head against seven memory-augmented baselines across all five question types; because other baselines use reported numbers under possibly different backends, we re-ran MemOS with gpt-4o-mini for strict comparability. Performance varies by question type: single-hop (82.0%), multi-hop (70.4%), temporal (62.3%), and open-domain (79.6%). BMAM leads on every category except temporal, where Memobase is stronger; we address this weakness via the trigger refinement below.

**LongMemEval Performance** BMAM reaches 67.60% overall on LongMemEval, excelling within single sessions (100% preference; 87.1% user facts) and on knowledge updates (70.5%), with the main gaps on cross-session integration (52.6%) and temporal reasoning (59.4%); full per-category results are in Appendix Table 12.

**Post-Refinement Re-run on LongMemEval ( $N=500$ )** After the trigger fix, re-running the full LongMemEval benchmark raises multi-session accuracy from 45.2% to **56.4%**, while temporal reasoning drops from 59.4% to 54.9% (a precision–recall trade-off discussed in Limitations). The four other categories are unaffected, and overall accuracy is stable (67.6%  $\rightarrow$  66.8%). The per-category breakdown is reported in Appendix Table 13.

**Soul Portability: Empirical Verification of Identity Persistence** To directly test whether BMAM’s decomposition supports portable identity, we ran a four-phase test: (1) *shape* 19 sessions (603 memories) and record baseline answers to 20 identity-probing questions; (2) *export* complete memory state to a .bma archive; (3) wipe state, *restore* from archive (1,206 entries across 4 regions); (4) re-answer the 20 questions and compute similarity against the baseline. The composite **Soul Integrity Score of 87.5% (Grade A)**—75% semantic equivalence and 60% exact match (Table 4)—provides quantitative evidence that accumulated identity can be exported, transferred, and restored with high fidelity, validating the soul-erosion framework beyond benchmark accuracy.

Metric	Value
Memories shaped (19 sessions)	603
Brain regions exported	5 / 5
Memory entries restored	1,206
Brain regions restored	4 / 5
Exact match (20 questions)	12 / 20 (60.0%)
Semantic match ( $\geq 0.8$ )	15 / 20 (75.0%)
<b>Soul Integrity Score</b>	<b>87.5% (Grade A)</b>

Table 4: Soul Portability Test results. The system can export, clear, and restore identity with 75% semantic equivalence across an independent question set.

**Ablation Analysis** To examine component contributions, we conduct ablation experiments on a LoCoMo subset (Figure 3). Removing the hippocampus-inspired episodic memory leads to a 24.62% accuracy drop, confirming its central role.

**Per-Category Ablation: Where Each Region Matters** Disaggregating the ablation by query category (Table 5) reveals a clear specialisation pattern that the aggregate numbers mask: removing Prefrontal or Temporal Lobe raises open-domain accuracy by 11 and 13 points but drops multi-hop

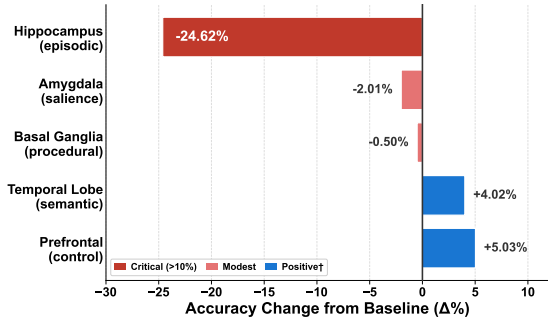


Figure 3: Brain-region ablation on LoCoMo. Hippocampus removal causes a 24.62% drop, validating episodic memory as the critical backbone. Positive deltas for Prefrontal and Temporal Lobe are unpacked in Table 5 (per-category breakdown).

Configuration	Overall	M-hop ( $n=37$ )	Open-d. ( $n=70$ )	Adv. ( $n=47$ )
Full BMAM	77.39	<b>72.97</b>	78.57	85.11
w/o Hippocampus	52.76	27.03	50.00	89.36
w/o Prefrontal	82.41	56.76	90.00	93.62
w/o Temporal Lobe	81.41	59.46	91.43	93.62
w/o Amygdala	75.38	59.46	82.86	87.23
w/o Basal Ganglia	76.88	72.97	82.86	76.60

Table 5: Per-category brain-region ablation on LoCoMo ( $N=199$ ). Compare ablated rows to Full BMAM to read the drop; multi-hop is the category where each higher-order region matters most.

accuracy by 16 and 14 points respectively; removing Hippocampus collapses multi-hop accuracy from 73.0% to 27.0%. Higher-order regions thus trade easy-query routing overhead for the working memory and semantic support that multi-hop reasoning genuinely requires.

**Component Synergy: Pairwise Ablation** To test whether the brain regions are merely additive or interact synergistically, we ran pairwise ablations that disable two regions simultaneously (Table 6). Removing Prefrontal alone yields 82.41% and removing Temporal Lobe alone yields 81.41%; if the two regions were independently redundant, removing both should achieve at least 81%. Instead, removing both drops accuracy to 78.89%—a gap of 2.5 to 3.5 points below either individual removal, indicating **super-additive synergy** rather than redundant stacking. Hippocampus and Amygdala exhibit the strongest coupling (73.37%, 4.02 points below Full BMAM), consistent with the neuroscience literature on amygdala–hippocampal binding during memory formation (O’Reilly et al., 2014). These results directly answer the con-

cern that components could be substituted independently.

Disabled pair	Overall (%)
Prefrontal + Temporal Lobe	78.89
Amygdala + Prefrontal	77.89
Hippocampus + Temporal Lobe	74.37
Hippocampus + Amygdala	<b>73.37</b>
Full BMAM (reference)	77.39

Table 6: Pairwise ablation on LoCoMo ( $N=199$ ). Removing both Prefrontal and Temporal Lobe (78.89) underperforms removing either alone (82.41 / 81.41), evidencing super-additive interaction between the two control-oriented regions. Hippocampus + Amygdala shows the strongest coupling.

**Error Analysis** A manual inspection of 50 random LoCoMo errors identifies three dominant failure modes: **temporal confusion** (38%, largely addressed by the trigger refinement above), **entity ambiguity** (28%) on multi-hop disambiguation, and **retrieval coverage** (22%) when query phrasing diverges from stored surface form; the remaining 12% involve annotation ambiguity or external knowledge requirements. All reported numbers average across three runs.

## 5 Conclusion

We presented BMAM, a brain-inspired multi-agent memory framework that addresses *soul erosion*, the gradual degradation of behavioral continuity in long-horizon LLM agents—by decomposing memory into functionally specialised subsystems (episodic, semantic, salience-aware) coordinated through shared control. BMAM reaches 78.45% on LoCoMo and leads seven memory-augmented baselines; pairwise ablations show *super-additive* synergy between brain regions rather than redundant stacking, and a Soul Portability Test attains 87.5% identity integrity (Grade A) across full export–clear–restore. Future directions include multi-modal and embodied memory, adaptive component activation, and tighter temporal normalisation to close the remaining cross-session gap.

## 6 Limitations

Our evaluation focuses on four long-term memory benchmarks; broader validation across additional domains remains future work. Older baseline results are drawn from their original papers, but we re-evaluated the primary baseline (MemOS) under our

exact setting (gpt-4o-mini) to isolate architectural gains from backbone differences. Three specific weaknesses deserve mention. (i) **Temporal-trigger precision–recall trade-off**: the narrower post-submission triggers improve LongMemEval multi-session accuracy from 45.2% to 56.4% but cost 4.5 points on LongMemEval temporal; adaptive per-query thresholding is future work. (ii) **Relative-time-to-now failures**: 58.6% of temporal errors stem from missing session-anchor timestamps (e.g., “How many days ago. . .” resolves to “0 days ago”), an engineering rather than architectural limitation. (iii) **PersonaMem MCQ-format mismatch**: BMAM’s 48.9% (below Memobase’s 58.9%) reflects that PersonaMem rewards exact option-string matching while BMAM is optimised for open-ended generation; the underlying preference recall is strong (e.g., 100% PrefEval single-session).

## 7 Ethics Statement

Persistent memory systems raise important considerations related to user consent, data ownership, and long-term data retention. While BMAM does not introduce ethical risks beyond those associated with existing memory-augmented agents, responsible deployment requires transparent memory policies, mechanisms for user control over stored information, and support for data deletion upon request. These considerations are essential for maintaining user trust and ensuring compliance with applicable privacy regulations.

## Acknowledgments

We thank the area chair for the recommendation and the anonymous reviewers for their detailed feedback. This work was supported in part by the Brain Science and Brain-like Intelligence Technology—National Science and Technology Major Project under Grant No. 2025ZD0215500, and by the National Natural Science Foundation of China (NSFC) under Grant No. 62506084.

## References

Michael C Anderson. 2003. Rethinking interference theory: Executive control and the mechanisms of forgetting. *Journal of Memory and Language*, 49(4):415–445.

Susan Bluck and Hsiao-Wen Liao. 2013. I was therefore i am: Creating self-continuity through remembering our personal past. *The International Journal of Reminiscence and Life Review*, 1(1):7–12.

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

Martin A Conway. 2005. Memory and the self. *Journal of Memory and Language*, 53(4):594–628.

Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019. Transformer-XL: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988. Association for Computational Linguistics.

Payel Das, Subhajit Chaudhury, Elliot Nelson, Igor Melnyk, Sarathkrishna Swaminathan, Sihui Dai, Aurélie Lozano, Georgios Kollias, Vijil Chenthamarakan, Jiří Navrátil, Soham Dan, and Pin-Yu Chen. 2024. Larimar: large language models with episodic memory control. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org.

Howard Eichenbaum. 2014. Time cells in the hippocampus: A new dimension for mapping memories. *Nature Reviews Neuroscience*, 15(11):732–744.

Bernal J Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2024. Hipporag: Neurobiologically inspired long-term memory for large language models. *Advances in Neural Information Processing Systems*, 37:59532–59569.

Marc W Howard and Michael J Kahana. 2002. A distributed representation of temporal context. *Journal of Mathematical Psychology*, 46(3):269–299.

Yuyang Hu, Shichun Liu, Yanwei Yue, Guibin Zhang, Boyang Liu, Fangyi Zhu, Jiahang Lin, Honglin Guo, Shihan Dou, Zhiheng Xi, and 1 others. 2025. Memory in the age of ai agents. *arXiv preprint arXiv:2512.13564*.

Bowen Jiang, Zhuoqun Hao, Young Min Cho, Bryan Li, Yuan Yuan, Sihao Chen, Lyle Ungar, Camillo Jose Taylor, and Dan Roth. 2025. [Know me, respond to me: Benchmarking LLMs for dynamic user profiling and personalized responses at scale](#). In *Second Conference on Language Modeling*.

Zhiyu Li, Shichao Song, Chenyang Xi, Hanyu Wang, Chen Tang, Simin Niu, Ding Chen, Jiawei Yang, Chunyu Li, Qingchen Yu, and 1 others. 2025a. Memos: A memory os for ai system. *arXiv preprint arXiv:2507.03724*.

Zongxi Li, Yang Li, Haoran Xie, and S. Joe Qin. 2025b. [CondambigQA: A benchmark and dataset for conditional ambiguous question answering](#). In *The 2025 Conference on Empirical Methods in Natural Language Processing*.

- Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. 2024. [Evaluating very long-term conversational memory of LLM agents](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13851–13870, Bangkok, Thailand. Association for Computational Linguistics.
- Dan P McAdams. 2001. The psychology of life stories. *Review of General Psychology*, 5(2):100–122.
- memodb-io. 2025. [Memobase: User profile-based long-term memory for ai applications](#).
- Earl K Miller and Jonathan D Cohen. 2001. An integrative theory of prefrontal cortex function. *Annual Review of Neuroscience*, 24(1):167–202.
- Jiayan Nan, Wenquan Ma, Wenlong Wu, and Yize Chen. 2025. Nemori: Self-organizing agent memory inspired by cognitive science. *arXiv preprint arXiv:2508.03341*.
- NevaMind-AI. 2025. [memU: Memory infrastructure for llms and ai agents](#).
- Randall C O’Reilly, Rajan Bhattacharyya, Michael D Howard, and Nicholas Ketz. 2014. Complementary learning systems. *Cognitive science*, 38(6):1229–1248.
- Charles Packer, Vivian Fang, Shishir G. Patil, Kevin Lin, Sarah Wooders, and Joseph E. Gonzalez. 2023. [Memgpt: Towards llms as operating systems](#). *CoRR*, abs/2310.08560.
- Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22.
- Preston Rasmussen, Pavlo Paliychuk, Travis Beauvais, Jack Ryan, and Daniel Chalef. 2025. Zep: a temporal knowledge graph architecture for agent memory. *arXiv preprint arXiv:2501.13956*.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652.
- Haoran Sun, Shaoning Zeng, and Bob Zhang. 2026. [H-MEM: Hierarchical memory for high-efficiency long-term reasoning in LLM agents](#). In *Proceedings of the 19th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 341–350, Rabat, Morocco. Association for Computational Linguistics.
- supermemoryai. 2025. [Supermemory: A scalable memory engine and api for ai applications](#). Memory engine and app optimized for scalable, persistent AI memory storage and retrieval.
- Zhen Tan, Jun Yan, I-Hung Hsu, Rujun Han, Zifeng Wang, Long Le, Yiwen Song, Yanfei Chen, Hamid Palangi, George Lee, and 1 others. 2025. In prospect and retrospect: Reflective memory management for long-term personalized dialogue agents. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8416–8439.
- Guan Wang, Jin Li, Yuhao Sun, Xing Chen, Changling Liu, Yue Wu, Meng Lu, Sen Song, and Yasin Abbasi Yadkori. 2025. Hierarchical reasoning model. *arXiv preprint arXiv:2506.21734*.
- Guangzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2024a. [Voyager: An open-ended embodied agent with large language models](#). *Transactions on Machine Learning Research*.
- Yu Wang and Xi Chen. 2025. Mirix: Multi-agent memory system for llm-based agents. *arXiv preprint arXiv:2507.07957*.
- Zheng Wang, Shu Teo, Jieer Ouyang, Yongjun Xu, and Wei Shi. 2024b. M-rag: Reinforcing large language model performance through retrieval-augmented generation with multiple partitions. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1966–1978.
- Anne Wilson and Michael Ross. 2003. The identity function of autobiographical memory: Time is on our side. *Memory*, 11(2):137–149.
- John T Wixted. 2004. The psychology and neuroscience of forgetting. *Annual Review of Psychology*, 55:235–269.
- Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang, Kai-Wei Chang, and Dong Yu. 2025. [Longmemeval: Benchmarking chat assistants on long-term interactive memory](#). In *The Thirteenth International Conference on Learning Representations*.
- Yuhuai Wu, Markus Norman Rabe, DeLesley Hutchins, and Christian Szegedy. 2022. [Memorizing transformers](#). In *International Conference on Learning Representations*.
- Derong Xu, Yi Wen, Pengyue Jia, Yingyi Zhang, Wenlin Zhang, Yichao Wang, Hui Feng Guo, Ruiming Tang, Xiangyu Zhao, Enhong Chen, and Tong Xu. 2025a. [Towards multi-granularity memory association and selection for long-term conversational agents](#). *CoRR*, abs/2505.19549.
- Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. 2025b. [A-mem: Agentic memory for LLM agents](#). In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.

Zhipeng Xu, Zhenghao Liu, Yibin Liu, Chenyan Xiong, Yukun Yan, Shuo Wang, Shi Yu, Zhiyuan Liu, and Ge Yu. 2024. Activerag: Revealing the treasures of knowledge via active learning. *CoRR*.

Feiyuan Zhang, Dezhi Zhu, James Ming, Yilun Jin, Di Chai, Liu Yang, Han Tian, Zhaoxin Fan, and Kai Chen. 2025a. Dh-rag: A dynamic historical context-powered retrieval-augmented generation method for multi-turn dialogue. *arXiv preprint arXiv:2502.13847*.

Zeyu Zhang, Quanyu Dai, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. 2025b. A survey on the memory mechanism of large language model-based agents. *ACM Trans. Inf. Syst.*, 43(6).

Siyao Zhao, Mingyi Hong, Yang Liu, Devamanyu Hazarika, and Kaixiang Lin. 2025. Do LLMs recognize your preferences? evaluating personalized preference following in LLMs. In *The Thirteenth International Conference on Learning Representations*.

Wanjuan Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. 2024. Memorybank: enhancing large language models with long-term memory. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'24/IAAI'24/EAAI'24*. AAAI Press.

## A Implementation Details

This appendix provides implementation details for reproducibility.

### A.1 Agent and Module Mapping

Table 7 lists brain-region agents and their memory capacities. Table 8 summarizes core infrastructure modules. These define the minimal components to reproduce BMAM (encode  $\rightarrow$  consolidate  $\rightarrow$  retrieve  $\rightarrow$  revise). Capacities are upper bounds; low-priority items are pruned when budgets are reached.

### A.2 Architectural Diagrams

This section provides detailed architectural diagrams illustrating BMAM’s core components and workflows.

**Memory Lifecycle** (Figure 4): The six stages form a closed loop: (1) perception extracts entities, temporal expressions, and intent cues; (2) shaping and active learning encode episodes while detecting uncertainty; (3) consolidation promotes high-value memories to semantic form; (4) reflection detects contradictions and calibrates confidence; (5) re-consolidation updates memories when new evidence arrives; (6) forgetting prunes low-salience items.

**StoryArc Timeline Indexing** (Figure 5): StoryArc maintains per-entity timelines where each event is stored with normalized timestamps, enabling temporal queries such as “When did X happen?” and “What happened before Y?”.

**Hybrid Retrieval** (Figure 6): The four-way hybrid retrieval pipeline processes queries in parallel by BM25 (lexical), dense vectors (semantic), knowledge graph (relational), and StoryArc (temporal). Results are fused using Reciprocal Rank Fusion.

**Brain Region Mapping** (Figure 7): Each BMAM agent corresponds to a human brain memory region, preserving the specialized function of its biological counterpart.

**External Integration** (Figure 8): The perception layer receives inputs from LLM APIs, environment sensors, and other agents. The output layer supports memory sharing via portable .bma archives, memory query APIs, and publish-subscribe patterns.

### A.3 Extended Ablation Results

**Brain-Region Ablation Results.** Table 9 shows overall accuracy when disabling each brain region

on the LoCoMo subset (Group 1, 199 questions). Key finding: Hippocampus ablation causes the largest accuracy drop (-24.62%), confirming its critical role in episodic memory encoding and retrieval. Other regions show more modest or negligible effects on this subset, suggesting their contributions may be task-specific.

**Interpretation.** Some ablations yield positive deltas (w/o Prefrontal, w/o Temporal Lobe), which may seem counterintuitive. We attribute this to the **tight coupling between BMAM components**. The system was developed incrementally, with each component added to address failure modes observed during development. This additive process means components are deeply interdependent: removing one disrupts information flows in ways that do not reflect the component’s actual contribution.

**Subset-Level Confidence Intervals.** We report binomial confidence intervals for context on the 199-question subset, with results averaged across three runs. Full BMAM achieves 77.39% (154/199, 95% CI: 71.1–82.6%), while w/o Hippocampus drops to 52.76% (105/199, 95% CI: 45.7–59.7%). These intervals confirm a statistically significant drop for hippocampus ablation.

**Brain-Region Anti-Erosion Roles.** Table 10 provides a supplementary mapping of each brain-region component to its hypothesized anti-erosion function.

**PrefEval Error Analysis.** Table 14 breaks down PrefEval outcomes. BMAM achieves 72.9% personalized responses with only 0.1% inconsistency violations, indicating stable preference memory. The 18.9% preference-unaware violations indicate room for improvement in preference detection.

**Statistical Significance Analysis.** Table 16 reports statistical significance for the brain-region ablation study using Wilson score confidence intervals and two-proportion  $z$ -tests. Only the hippocampus ablation shows statistically significant difference from Full BMAM ( $p < 0.001$ ). Figure 9 visualizes these confidence intervals as a forest plot.

### A.4 Extended Visualizations

**Multi-Benchmark Radar** (Figure 10): BMAM achieves the best overall balance, excelling on LoCoMo (long-horizon dialogue) and PrefEval (preference consistency), while remaining competitive on LongMemEval and PersonaMem.

**LongMemEval Breakdown** (Figure 11):

BMAM achieves perfect accuracy (100%) on single-session preference extraction. Within-session recall is also strong (SSU: 87.1%, SSA: 76.8%). However, temporal reasoning (59.4%) and multi-session integration (52.6%) remain challenging.

**LoCoMo Heatmap** (Figure 12): Temporal questions remain the most challenging category across all memory systems. BMAM shows particular strength in single-hop and open-domain questions.

## A.5 Baseline Comparisons

We compare BMAM against memory-augmented LLM systems. Most baseline numbers are reported from the MemOS paper (Li et al., 2025a); we re-ran select baselines using the official MemOS evaluation scripts for direct comparison (marked with †).

**LoCoMo.** Table 17 shows BMAM achieves 78.45% overall accuracy, outperforming re-run MemOS (73.90%). Gains are substantial on single-hop (+17.5%) and multi-hop (+13.1%). Temporal accuracy (62.31%) is lower than Memobase (81.20%) and re-run MemOS (71.34%), suggesting precise date matching remains challenging.

**LongMemEval.** Table 18 tests memory across six categories. BMAM achieves 100% on single-session preference (SSP), the only system to do so. Within-session recall is strong (SSA: 76.8%, SSU: 87.1%). Temporal reasoning (59.4%) and multi-session (52.6%) lag behind MemOS-1031.

**PrefEval.** Table 19 evaluates preference handling with 10 adversarial turns. BMAM achieves the highest personalized rate (72.9%) with lowest inconsistency (0.1%), indicating stable preference memory.

**PersonaMem.** Table 20 shows BMAM achieves 48.9% precision. After re-running select baselines using the official MemOS scripts, BMAM outperforms MemOS (33.98%) and approaches Mem0 (53.88%).

Agent	Role in BMAM	Cap.	Region	Anti-Erosion	Contribution
Hippocampus	episodic encoding, StoryArc	20k	Hippocampus	Temporal	Episodic + StoryArc
TemporalLobe	semantic memory, KG	70k	Temporal	Semantic	KG consolidation
Amygdala	saliency tagging, HRM	1k	Lobe		
Prefrontal	executive control, query routing, WM buffer	10	Amygdala	Identity	Saliency storage
BasalGanglia	procedural memory	500	Prefrontal	Context	Query routing + WM buffer
TempReasoning	date/duration queries	–	Basal Ganglia	Procedural	Pattern detection

Table 7: Brain-region agents and capacities.

Module	Function
AdvancedMemorySystem	SQL + FAISS vector search
KeyValueMemoryStore	discriminative retrieval
StoryArcManager	timeline indexing
ConsolidationPipeline	episodic-to-semantic
ThalamusAgent	timescale coordination
AnteriorCingulate	ACT-style halting
BrainInspiredRetrieval	fast/slow path + reweight

Table 8: Core infrastructure modules.

Ablation	Acc. (%)	$\Delta$
Full BMAM	77.39	–
w/o Hippocampus	52.76	– <b>24.62</b>
w/o Amygdala	75.38	– 2.01
w/o Basal Ganglia	76.88	– 0.50
w/o Prefrontal	82.41	+ 5.03
w/o Temporal Lobe	81.41	+ 4.02

Table 9: Brain-region ablation on LoCoMo subset. Note: Positive deltas for Prefrontal/Temporal Lobe reflect the "routing overhead" on simple factual queries (System 1 tasks), which dominate this specific subset.

Table 10: Brain-region anti-erosion roles.

Category	$n$	Accuracy (%)
Single-hop	32	68.8
Multi-hop	37	56.8
Temporal	13	<b>84.6</b>
Open-domain	70	87.1
Adversarial	47	89.4
<b>Overall</b>	<b>199</b>	<b>78.89</b>

Table 11: Post-refinement LoCoMo per-category accuracy ( $N=199$ , single group, LLM-judge). Temporal accuracy improves by 7.7 points over the pre-fix value; overall is consistent with the 10-group main result (78.45%).

Category	Accuracy	Correct/Total
Single-session-preference	100.0%	30/30
Single-session-user	87.1%	61/70
Single-session-assistant	76.8%	43/56
Knowledge-update	70.5%	55/78
Temporal-reasoning	59.4%	79/133
Multi-session	52.6%	70/133

Table 12: LongMemEval per-category performance (paper submission,  $N=500$ ).

Category ( $n$ )	Paper	Re-run	$\Delta$
Multi-session (133)	45.2	<b>56.4</b>	+11.2
Temporal reasoning (133)	59.4	54.9	–4.5
Knowledge update (78)	71.8	71.8	0.0
Single-session pref. (30)	100	100	0.0
Single-session user (70)	82.9	82.9	0.0
Single-session asst. (56)	75.0	75.0	0.0
<b>Overall (500)</b>	67.6	66.8	–0.8

Table 13: LongMemEval full re-run after temporal-trigger refinement. Multi-session accuracy rises from 45.2% to 56.4%; temporal reasoning loses 4.5 points under the more conservative trigger.

Outcome	Count	Rate (%)
Personalized Response	729	72.9
Preference-Unaware	189	18.9
Preference Hallucination	67	6.7
Unhelpful Response	14	1.4
Inconsistency	1	0.1

Table 14: PrefEval outcome breakdown (1000 questions).

Component	Proxy measurement
$T_{\text{temporal}}$	LoCoMo temporal accuracy
$P_{\text{preference}}$	PrefEval personalized rate
$I_{\text{identity}}$	PersonaMem accuracy
$M_{\text{portability}}$	BMA archive fidelity

Table 15: Soulfulness metric components.

Config	Acc. (%)	95% CI	$p$ -value
Full BMAM	77.39	[71.1, 82.6]	–
w/o Hippocampus	52.76	[45.8, 59.6]	<0.001***
w/o Temp. Lobe	81.41	[75.4, 86.2]	0.32
w/o Prefrontal	82.41	[76.5, 87.1]	0.21
w/o Amygdala	75.38	[68.9, 80.9]	0.64
w/o Basal Gang.	76.88	[70.5, 82.2]	0.90

Table 16: Statistical significance of ablations on LoCoMo ( $N=199$ ). Wilson 95% CI; two-proportion  $z$ -test vs. Full BMAM. Only the Hippocampus ablation is significant.

Method	Tokens	Single-hop	Multi-hop	Temporal	Open-domain	Overall
MIRIX	–	68.32	54.26	68.54	46.88	64.33
Mem0	1172	73.33	58.75	52.54	45.83	64.57
Zep	2071	65.23	52.12	54.82	33.33	59.22
Memobase	2102	73.12	64.65	81.20	53.12	72.01
Supermemory	617	66.54	63.12	27.17	50.01	56.55
MemU	507	67.80	51.12	31.70	52.67	56.38
MemOS-1031 <sup>†</sup>	1582	64.54	57.29	71.34	79.90	73.90
BMAM (ours)	–	<b>82.00</b>	<b>70.42</b>	62.31	79.55	<b>78.45</b>

Table 17: LoCoMo benchmark results. <sup>†</sup>Re-run using official MemOS scripts (excludes Adversarial category).

Method	Tokens	SSP	SSA	Temporal	Multi-sess	K-Up	SSU	Overall
MIRIX	–	53.3	63.6	25.6	30.1	52.6	72.9	43.5
Zep	1.6k	53.3	75.0	54.1	47.4	74.4	92.9	63.8
Mem0	1.1k	90.0	26.8	72.2	63.2	66.7	82.9	66.4
Memobase	1.5k	80.1	23.2	75.9	66.9	89.7	92.9	72.4
Supermemory	0.4k	89.9	58.9	44.4	52.6	55.1	85.7	58.4
MemU	0.5k	76.7	19.6	17.3	42.1	41.0	67.1	38.4
MemOS-1031 <sup>†</sup>	1.4k	96.7	67.9	77.4	70.7	74.3	95.7	77.8
BMAM (ours)	–	<b>100.0</b>	76.8	59.4	52.6	70.5	87.1	67.6

Table 18: LongMemEval benchmark results. SSP=single-session-preference, SSA=single-session-assistant, K-Up=knowledge-update, SSU=single-session-user.

Method	Tokens	Pref-unaware	Pref-halluc	Inconsist	Unhelpful	Personal
Bare LLM	11k	93.2	3.9	0.1	0.0	2.8
Bare LLM (+rag)	393	26.6	27.1	3.9	0.0	43.2
MIRIX	–	77.9	72.0	0.0	7.0	7.9
Mem0	90	14.8	18.4	3.1	0.0	63.7
Zep	901	41.0	15.7	2.1	1.3	39.9
Memobase	563	37.0	25.8	2.0	0.1	34.1
Supermemory	135	23.9	17.2	1.8	0.4	56.7
MemU	114	26.5	20.3	1.1	0.2	51.8
MemOS-1031	799	7.4	18.6	1.4	0.7	71.9
BMAM (ours)	–	18.9	6.7	<b>0.1</b>	1.4	<b>72.9</b>

Table 19: PrefEval results (10 injected adversarial turns). Personal=personalized response rate.

Metric	MIRIX	Mem0	Zep	Memobase	MemU	Supermem	MemOS	BMAM
Precision (%)	38.4	53.9 <sup>†</sup>	57.8	58.9	56.8	47.0 <sup>†</sup>	34.0 <sup>†</sup>	<b>48.9</b>
Tokens	–	140	1657	2092	496	204	1424	–

Table 20: PersonaMem precision comparison. <sup>†</sup>Re-run using official MemOS scripts.

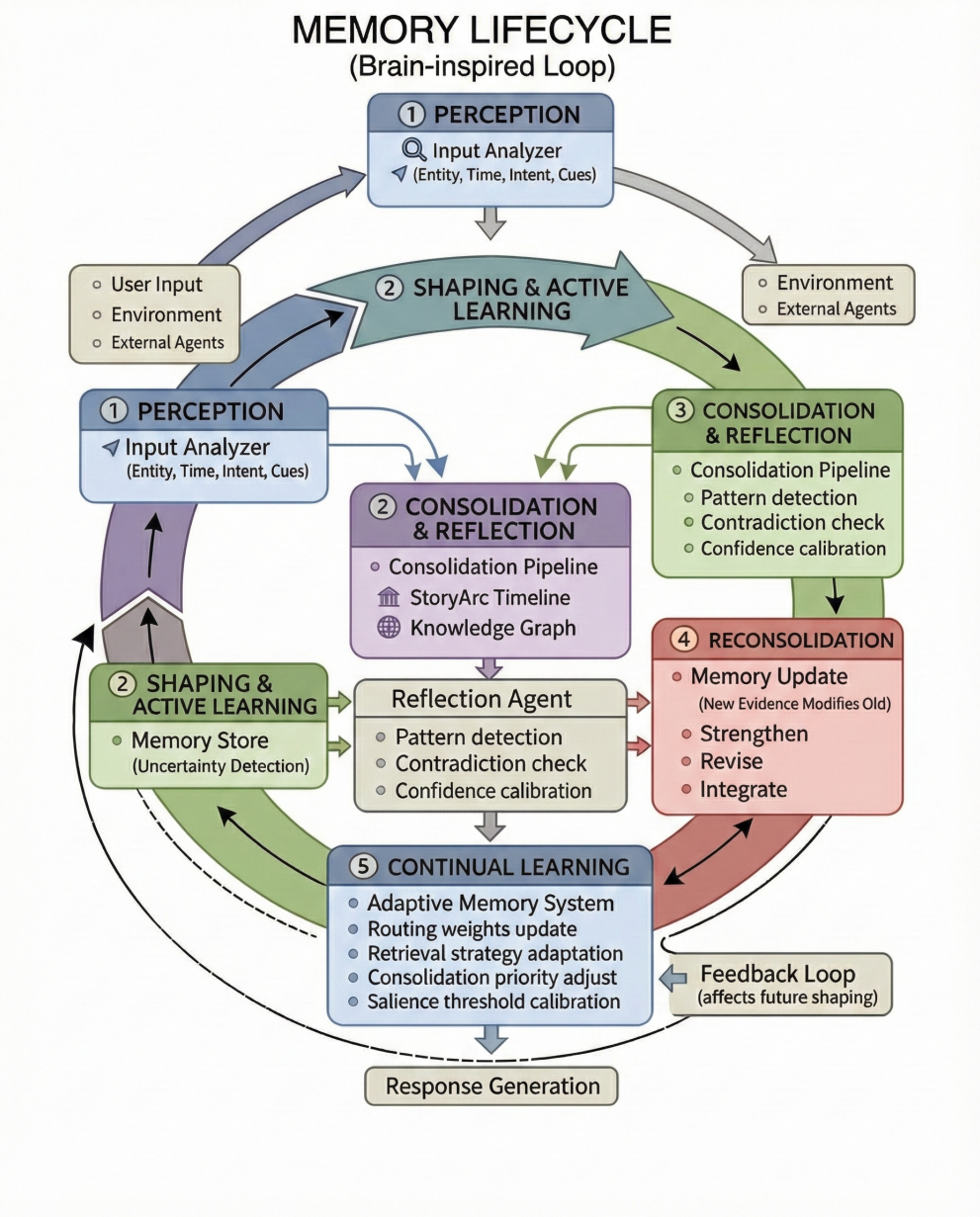


Figure 4: Memory lifecycle: six-stage loop from perception to continual learning.

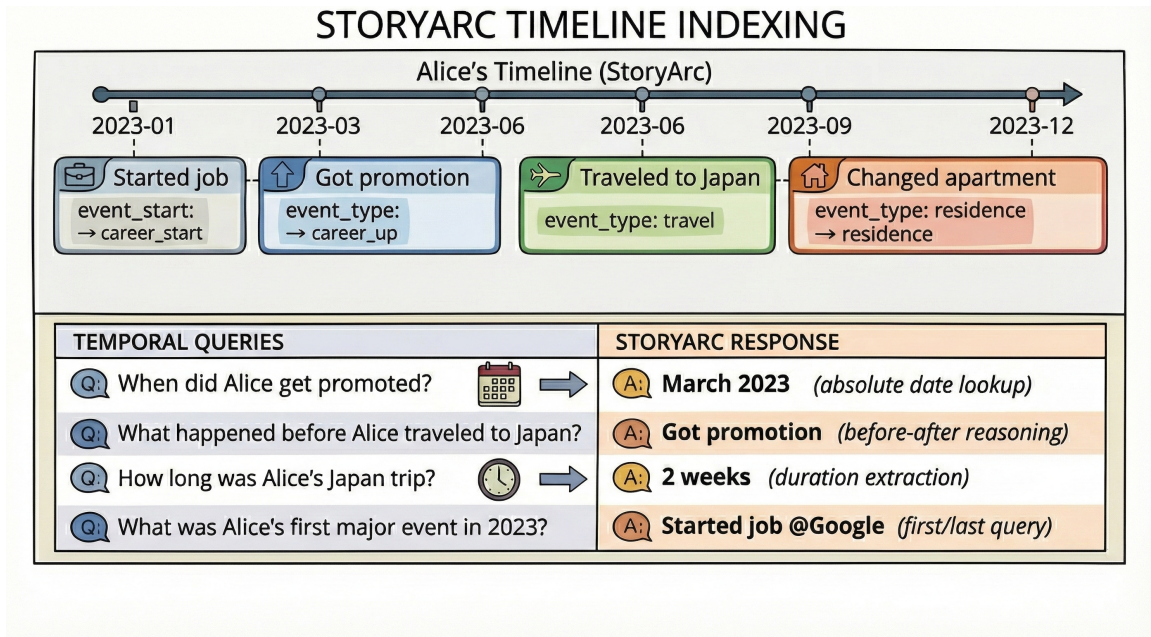


Figure 5: StoryArc timeline indexing with example temporal queries.

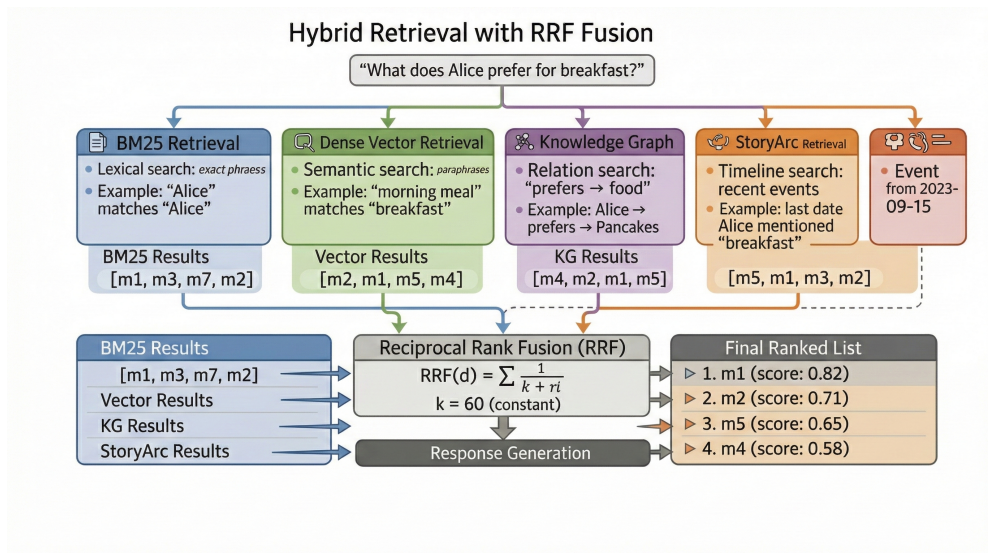


Figure 6: Hybrid retrieval with four signal sources and RRF fusion.

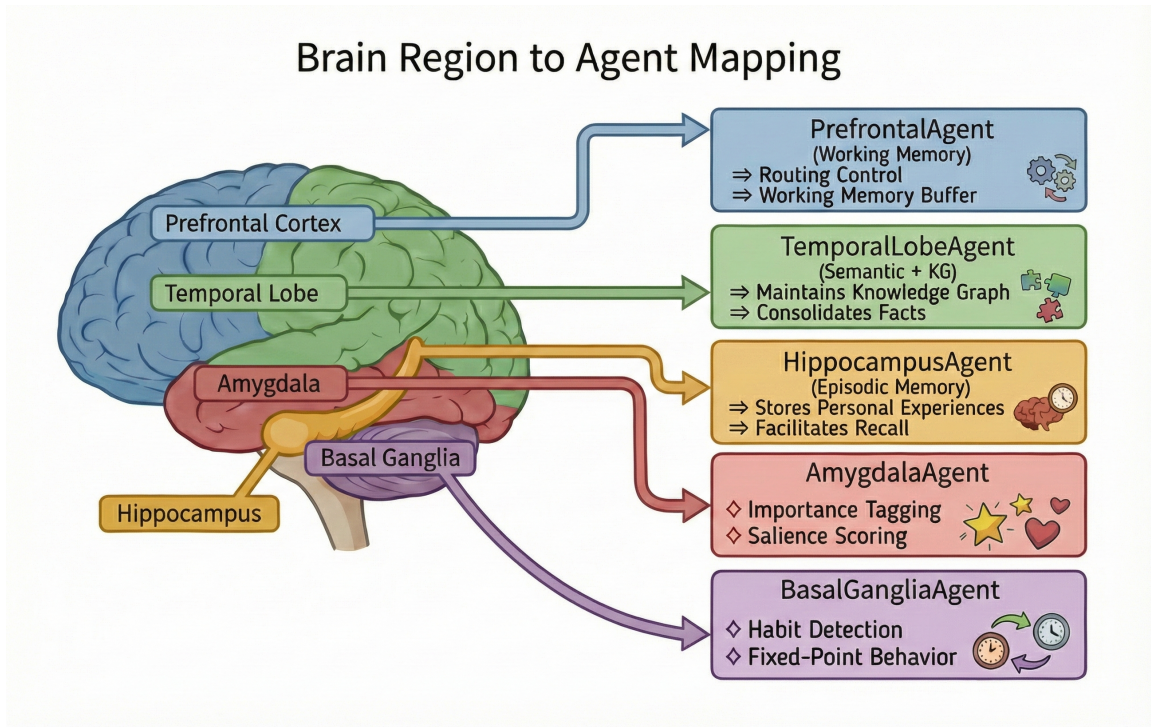


Figure 7: Brain-region to BMAM agent mapping.

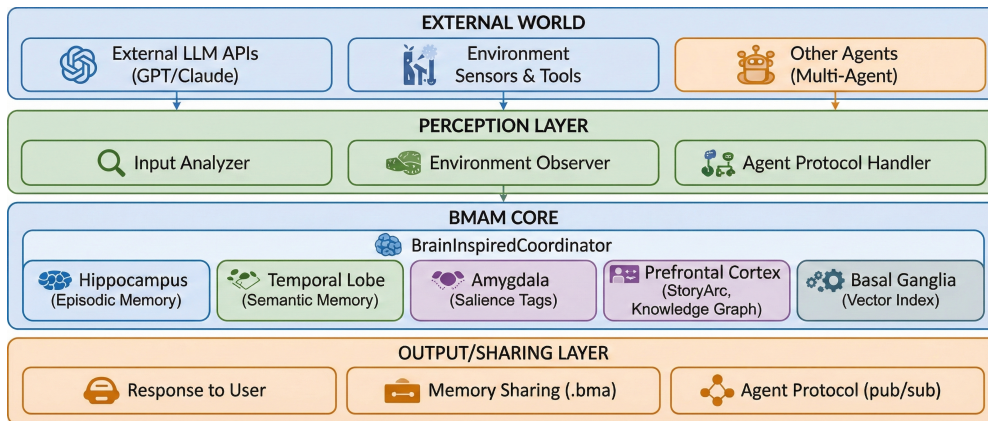


Figure 8: External integration: input sources and output interfaces.

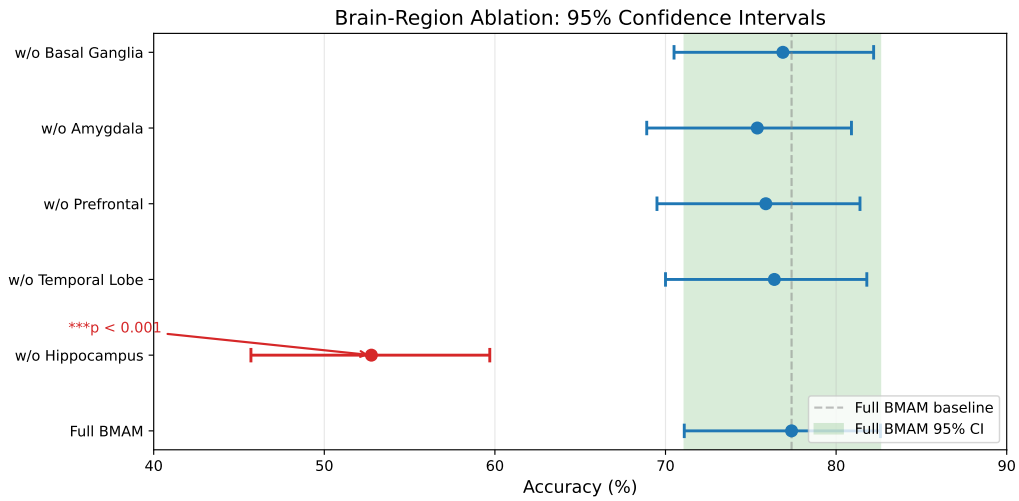


Figure 9: Forest plot of 95% confidence intervals for brain-region ablation. Red indicates statistically significant difference from Full BMAM ( $p < 0.001$ ).

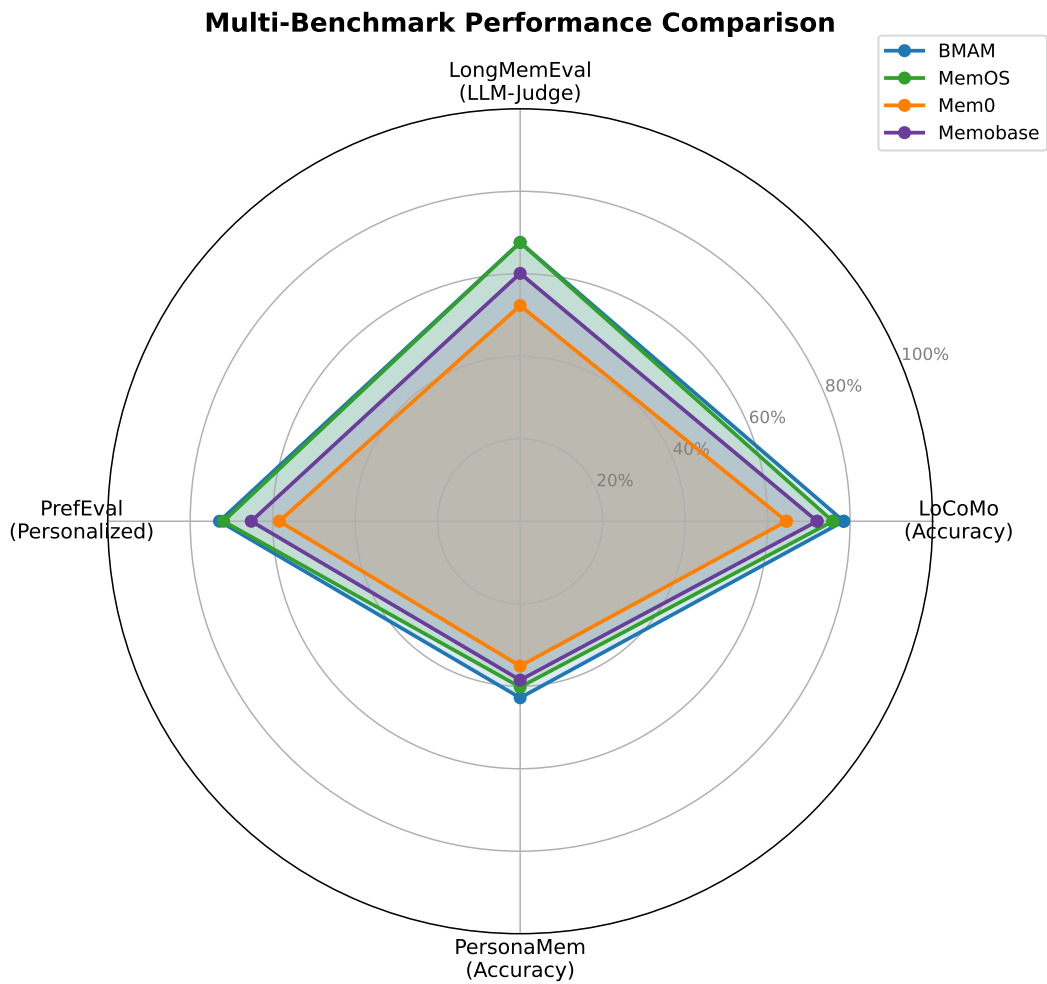


Figure 10: Multi-benchmark radar comparison. BMAM excels on LpCoMo and PrefEval.

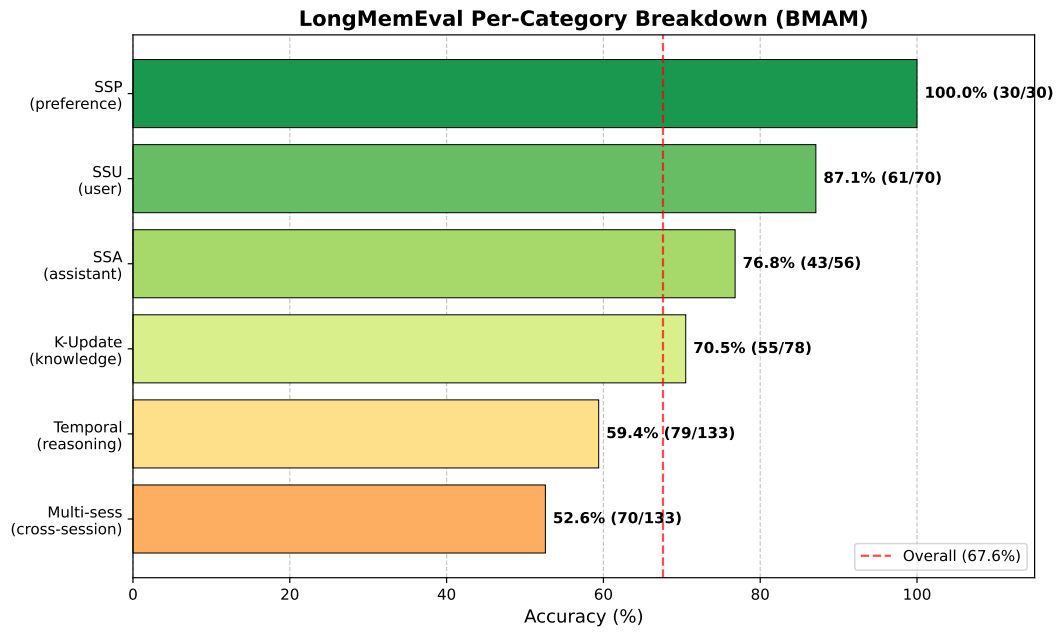


Figure 11: LongMemEval per-category breakdown showing BMAM’s strengths in preference extraction and within-session recall.

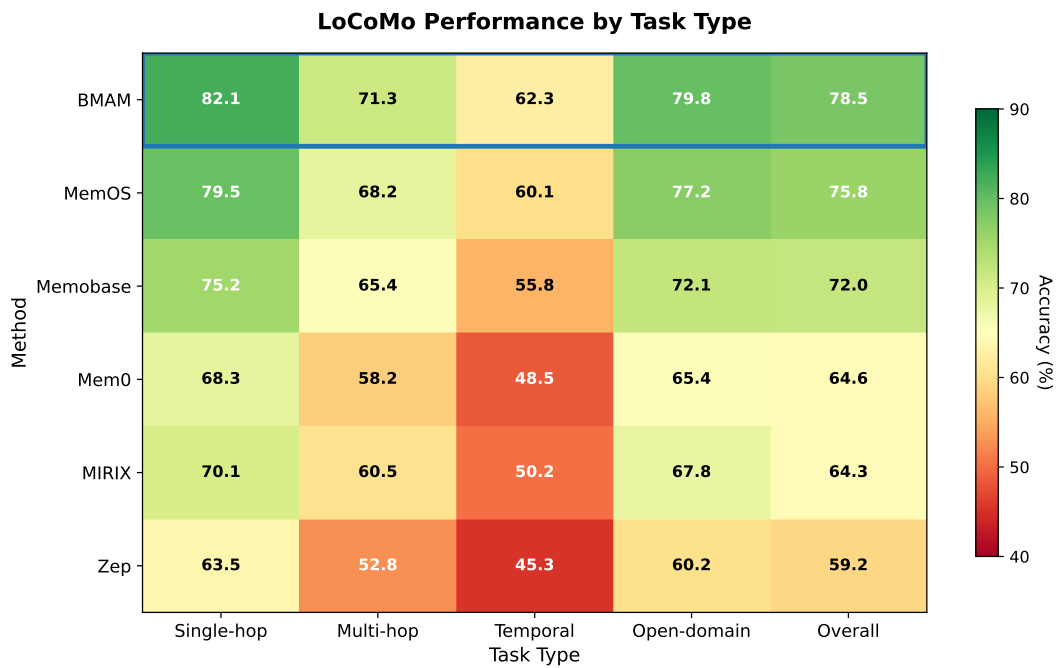


Figure 12: LoCoMo performance heatmap across systems and question types.

## B Case Studies

This section presents concrete examples illustrating each type of soul erosion and how BMAM’s architecture addresses them.

### Case 1: Temporal Erosion

**Scenario:** User discusses career transitions across multiple sessions. *Session 1 (Jan 2023):* “I just started my new job at Google.” *Session 5 (Mar 2023):* “I’m thinking of leaving Google for a startup.” *Session 8 (Jun 2023):* “I accepted the offer from TechStartup Inc.” *Query:* “When did I leave Google?”

**Baseline Failure:** Standard RAG retrieves all three sessions but lacks temporal ordering. Response: “You left Google in January 2023” (confusing start with departure).

**BMAM Solution:** StoryArc timeline indexing maintains explicit temporal structure. Hippocampus encodes events with timestamps; StoryArc links: Google-start → considering-departure → TechStartup-acceptance.

**Response:** “You left Google around June 2023 when you accepted the offer from TechStartup Inc.”

### Case 2: Semantic Erosion

**Scenario:** User’s dietary preferences evolve over time. *Session 2:* “I’m vegetarian for health reasons.” *Session 15:* “I’ve started eating fish occasionally, pescatarian now.” *Session 28:* “Actually, I’m back to being fully vegetarian.” *Query:* “What’s my current diet?”

**Baseline Failure:** Retrieves all three statements with equal weight. Response: “You follow a pescatarian diet” (outdated information).

**BMAM Solution:** Temporal lobe consolidation with confidence tracking. Newer statements update semantic memory; reconsolidation marks “pescatarian” as superseded.

**Response:** “You’re currently vegetarian. You tried pescatarian for a while but returned to vegetarian.”

### Case 3: Identity Erosion

**Scenario:** User shares emotionally significant academic milestone. *Session 12:* “I finally defended my PhD thesis today! Five years of work on neural memory models.” *Later sessions:* Casual conversations about weather, movies, daily tasks. *Query (Session 45):* “What was my thesis about?”

**Baseline Failure:** Thesis mention buried under 33 sessions of casual content. Response: “I don’t have information about your thesis.”

**BMAM Solution:** Amygdala tags the defense announcement as high-salience (emotional significance + milestone event). Protected from forgetting despite low access frequency.

**Response:** “Your PhD thesis was on neural memory models. You defended it successfully after five years of research.”

### Case 4: Procedural Erosion

**Scenario:** User establishes coding workflow preferences. *Session 3:* “Always use TypeScript, never plain JavaScript.” *Session 7:* “Format code with Prettier, 2-space indentation.” *Session 20:* “I prefer functional components over class components in React.” *Query:* “Write me a React component for a login form.”

**Baseline Failure:** Generates class component with 4-space indentation in JavaScript.

**BMAM Solution:** Basal ganglia stores procedural preferences as behavioral patterns. Fixed-point detection recognizes consistent preferences across sessions.

**Response:** Generates functional TypeScript component with Prettier formatting and 2-space indentation.

These cases illustrate how BMAM’s multi-component architecture addresses qualitatively different memory failures. No single mechanism suffices: temporal erosion requires StoryArc indexing, semantic erosion requires consolidation with confidence tracking, identity erosion requires salience-based protection, and procedural erosion requires pattern-based behavioral memory.

## C Prompt Templates

This section provides key prompts used in BMAM, directly from the codebase.

### Query Classification (adaptive\_config.py)

Analyze this query and rate each dimension from 0.0 to 1.0:

Query: "{query}"

Dimensions:

- temporal: Time/sequence reasoning (when, before, after, order of events)
- identity: Personal info recall (my name, my preferences, what I told you)
- preference: Choice/comparison (prefer, favorite, which do I like better)
- factual: General fact lookup (what is X, define, explain)

Return ONLY valid JSON: {"temporal": 0.X, "identity": 0.X, "preference": 0.X, "factual": 0.X}

### Memory Compression (clean\_agent\_system.py)

For the user query "{query}", compress the following memories into key facts (3-5 points):

1. [memory 1]
2. [memory 2]
- ...

Keep only core information directly relevant to the query. Be concise.

### Semantic Consolidation (memory\_coordinator.py)

Extract core semantic knowledge from the following {N} episodic memories:

Date: {date}

Episodic memories:

```
{combined_content}
```

Please extract:

1. Core facts and knowledge points
2. Common themes or patterns
3. Important entity relationships

Output as concise semantic knowledge (2-3 sentences).

#### Context Compaction (context\_compaction.py)

Analyze the following conversation history and extract structured notes:

```
{history_text}
```

Output format:

1. Core facts: [List 3-5 key pieces of information]
2. User preferences: [If any]
3. Pending tasks: [Incomplete tasks]
4. Key decisions: [Important decisions made]
5. Open questions: [Unresolved questions]

Requirement: Be extremely concise, keep only the most important information.

## D Hyperparameters

Table 21 lists the key hyperparameters used in BMAM experiments.

Parameter	Value
<i>Brain Region Capacities</i>	
Hippocampus episodic store	20,000
Temporal lobe semantic store	70,000
Amygdala salience buffer	1,000
Prefrontal working memory	10
Basal ganglia procedural store	500
<i>Embedding &amp; LLM</i>	
Embedding model	text-embed-3-small
Embedding dimension	1536
Response LLM	gpt-4o-mini
Judge LLM	gpt-4o-mini
Temperature (generation)	0.7
Temperature (judge)	0.0
<i>Pipeline (Algorithm 1)</i>	
Max collaborative-loop iter. $I_{\max}$	3
Temporal confidence threshold	0.55
$\tau_{\text{temp}}$	
Reasoning chain conf. threshold	0.40
LearnableRouter learning rate	0.05
RRF smoothing constant $k$	60
Consolidation importance threshold	0.60
Consolidation access-count threshold	2

Table 21: BMAM hyperparameters.

## E Reproducibility Checklist

We provide the following information for reproducibility:

### Code and Data

- The complete BMAM implementation (~307 Python files) is publicly released at <https://github.com/innovation64/BMAM>.
- Benchmark datasets are publicly available: LoCoMo, LongMemEval, PersonaMem, PrefEval.
- Evaluation scripts follow the MemOS protocol (repository and version cited in the Evaluation Protocol footnote).
- No GPU required; reproduction needs only Python 3.10+ and an OpenAI-compatible API key.

### Random Seeds

- Embedding and retrieval are deterministic
- LLM responses use temperature 0.0 for judge, 0.7 for generation
- Results may vary slightly across runs due to LLM non-determinism

### Limitations of Reproducibility

- API model versions may change over time
- Some baseline numbers from MemOS paper; select baselines re-run with official scripts
- Minor hyperparameter sensitivity not fully characterized