

# HiMem: Hierarchical Long-Term Memory for LLM Long-Horizon Agents

Anonymous ACL submission

## Abstract

Although long-term memory systems have made substantial progress in recent years, they still exhibit clear limitations in adaptability, scalability, and self-evolution under continuous interaction settings. Inspired by cognitive theories, we propose HiMem, a hierarchical long-term memory framework for long-horizon dialogues, designed to support memory construction, retrieval, and dynamic updating during sustained interactions. HiMem constructs cognitively consistent Episode Memory via a Topic-Aware Event-Surprise Dual-Channel Segmentation strategy, and builds Note Memory that captures stable knowledge through a multi-stage information extraction pipeline. These two memory types are semantically linked to form a hierarchical structure that bridges concrete interaction events and abstract knowledge, enabling efficient retrieval without sacrificing information fidelity. HiMem supports both hybrid and best-effort retrieval strategies to balance accuracy and efficiency, and incorporates conflict-aware Memory Reconsolidation to revise and supplement stored knowledge based on retrieval feedback. This design enables continual memory self-evolution over long-term use. Experimental results on long-horizon dialogue benchmarks demonstrate that HiMem consistently outperforms representative baselines in accuracy, consistency, and long-term reasoning, while maintaining favourable efficiency. Overall, HiMem provides a principled and scalable design paradigm for building adaptive and self-evolving LLM-based conversational agents.

## 1 Introduction

Large language models (LLMs) have demonstrated remarkable progress in language understanding and reasoning, enabling the development of LLM-based agents for complex, multi-turn tasks such as personalized assistance, planning, and long-term decision support (Yu et al., 2025; Su et al., 2025;

Liu et al., 2025). In realistic interactive settings, however, these agents are required to operate over extended time horizons, where relevant information is scattered across long dialogues and multiple sessions. Despite strong short-term reasoning ability, existing LLM agents still struggle to reliably preserve, organize, and utilize information over long time spans. This limitation has emerged as a fundamental bottleneck for building adaptive and consistent long-horizon conversational agents (Xu et al., 2022; Wu et al., 2025a; Li et al., 2025).

Recent efforts to address this challenge can be broadly categorized into three directions. Retrieval-augmented generation (RAG) systems introduce external memory stores to fetch relevant information on demand, improving factual grounding (Alonso et al., 2024; Jiménez Gutiérrez et al., 2024; Sarthi et al., 2024; Edge et al., 2024). Long-context modeling approaches extend the context window to thousands or even millions of tokens, enabling direct reasoning over extended histories (Du et al., 2025; Fountas et al., 2025; Qian et al., 2025; Lee et al., 2024). More recently, structured long-term memory systems have been proposed to persistently store and retrieve dialogue information in compressed or structured forms (Hatalis et al., 2023; Nan et al., 2025). While these methods significantly improve efficiency and continuity, they still exhibit systematic limitations when applied to long-horizon interactions (Yue et al., 2024; Zhang et al., 2025c).

From both empirical observations and cognitive perspectives, we identify three recurring challenges that existing long-term memory systems struggle to address simultaneously. First, **semantic misalignment** arises when extracted memories are detached from their original dialogue context, leading to errors in resolving temporal references, coreference, and implicit semantics. Second, most systems rely on **monolithic or insufficiently hierarchical memory structures**, forcing a trade-off be-

tween information fidelity and retrieval efficiency. Fine-grained dialogue logs preserve rich context but incur high retrieval costs, whereas aggressively abstracted representations reduce cost at the expense of critical details needed for reasoning and personalization. Third, memory updates are typically **static or similarity-driven**, lacking principled mechanisms to revise or correct stored knowledge when new information partially overlaps with, extends, or contradicts existing memories. As a result, long-term consistency degrades over sustained interactions.

Inspired by cognitive theories of human memory (Gilboa and Marlatte, 2017; Ghosh and Gilboa, 2014; Fauconnier and Turner, 2003; Nader, 2015), we argue that effective long-term memory for LLM agents must satisfy three properties: (i) a *hierarchical structure* that bridges concrete interaction events and abstracted knowledge, (ii) a *unified semantic alignment mechanism* that preserves interpretability across memory representations, and (iii) a *conflict-aware update process* that supports continual self-evolution rather than static accumulation. Based on these principles, we propose **HiMem**, a hierarchical long-term memory framework designed for long-horizon conversational agents.

HiMem organizes memory into two semantically linked layers. *Episode Memory* preserves fine-grained, temporally grounded interaction segments constructed via a Topic-Aware Event–Surprise Dual-Channel Segmentation strategy, which aligns memory boundaries with both topical shifts and cognitively salient discontinuities. *Note Memory* abstracts stable knowledge such as facts, user preferences, and user profiles through a multi-stage information extraction pipeline. These two memory types form a hierarchical transition from concrete events to compact knowledge representations, enabling efficient retrieval without sacrificing information fidelity. During retrieval, HiMem supports both a hybrid retrieval strategy and a best-effort retrieval strategy that descends from abstract knowledge to concrete events only when necessary. Crucially, retrieval failures are treated as learning signals: HiMem performs conflict-aware Memory Reconsolidation to supplement missing knowledge and revise existing memories, enabling continuous self-evolution over long-term use.

We evaluate HiMem on long-horizon dialogue benchmarks and demonstrate that it consistently outperforms representative baselines in accuracy, consistency, and efficiency. Extensive ablation

studies further validate the necessity of hierarchical memory organization, semantic alignment, and conflict-aware updating for robust long-term reasoning.

In summary, this paper makes the following contributions:

- We propose **HiMem**, a hierarchical long-term memory framework that integrates episodic and knowledge-oriented memories to support scalable and adaptive long-horizon conversational agents.
- We introduce a Topic-Aware Event–Surprise Dual-Channel Segmentation mechanism and a multi-stage information extraction pipeline to construct cognitively consistent and efficient memory representations.
- We design a conflict-aware Memory Reconsolidation mechanism that enables long-term memory systems to self-correct and evolve during sustained interactions.
- We provide extensive experimental evidence showing that principled hierarchical design and dynamic updating substantially improve long-horizon reasoning performance.

## 2 Methodology

HiMem is a modular long-term memory framework built upon a hierarchical architecture that integrates episodic interaction records with abstracted knowledge representations. It is designed to support efficient retrieval, semantic consistency, and continual memory evolution during long-horizon interactions.

### 2.1 Overall Framework

As shown in Figure 1, HiMem consists of three core modules: (i) a hierarchical memory construction module that builds Episode Memory and Note Memory from raw dialogues, (ii) a hierarchical memory retrieval module that supports both hybrid and best-effort retrieval strategies, and (iii) a conflict-aware memory updating module that enables continual self-evolution. Episode Memory preserves fine-grained interaction events, while Note Memory consolidates stable knowledge such as facts, user preferences, and user profiles. The two memory layers are semantically linked to form a hierarchy that bridges concrete experiences and abstract knowledge.

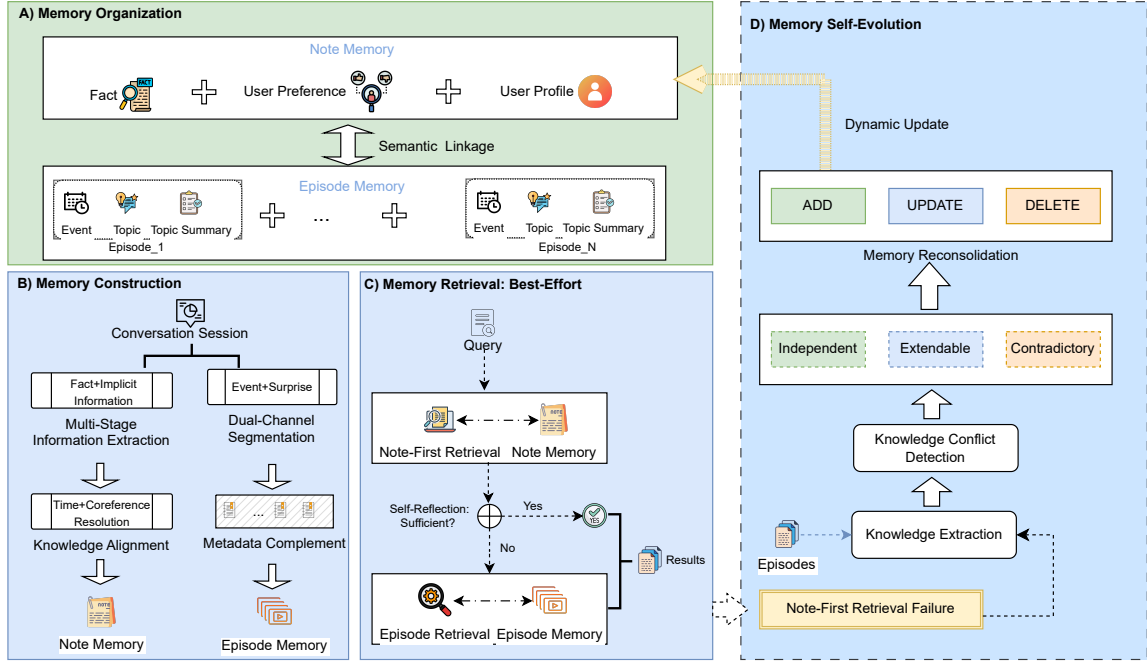


Figure 1: **Overview of HiMem.** (A) *Memory organization*: a hierarchical connection between Episode Memory and Note Memory. (B) *Memory construction*: pipelines that transform dialogue logs into Episode Memory and Note Memory. (C) *Best-effort retrieval*: hierarchical retrieval in the order of **Note Memory** → **Episode Memory**, with an LLM assessing evidence sufficiency. (D) *Memory self-evolution*: when evidence from Note Memory is insufficient, the system supplements potentially missing information from Episode Memory and triggers conflict detection and updating.

## 2.2 Memory Construction

Memory construction in HiMem follows a multi-stage pipeline that progressively transforms raw dialogue logs into structured long-term memory representations. This pipeline unifies event-level segmentation, knowledge extraction, and semantic alignment to ensure both fidelity and efficiency.

### 2.2.1 Episode Memory

Episode Memory records fine-grained interaction events aligned with topical and cognitive boundaries. Given a dialogue sequence, HiMem segments it into a sequence of non-overlapping episodes. Each episode is represented by a structured record containing an ID, timestamp, topic, topic summary, metadata, and the corresponding dialogue segment.

**Dual-Channel Segmentation.** To obtain cognitively coherent episodes, HiMem adopts a Topic-Aware Event–Surprise Dual-Channel Segmentation strategy. A segmentation boundary is introduced when either (i) a topical shift occurs in discourse goals or subtopics, or (ii) a salient discontinuity is detected, such as an abrupt change in intent or

emotional state. These two criteria are fused using an OR rule, producing event units that align with both semantic continuity and cognitive salience.

Segmentation is performed in a single pass, where an LLM jointly evaluates topical and surprise signals and directly outputs the final segmentation. This design yields compact and self-contained episodes that reduce cross-segment interference while preserving critical contextual evidence for downstream reasoning.

### 2.2.2 Note Memory

Note Memory focuses on long-term storage of knowledge-oriented information that remains stable or reusable across interactions. From each dialogue, HiMem extracts three categories of knowledge:

$$K = \{K_{\text{fact}}, K_{\text{pref}}, K_{\text{profile}}\},$$

where  $K_{\text{fact}}$  denotes objective facts and events,  $K_{\text{pref}}$  captures user preferences, and  $K_{\text{profile}}$  represents relatively stable user traits.

**Multi-Stage Knowledge Extraction.** Knowledge extraction is decomposed into three stages to avoid semantic collapse. Stage 1 extracts in-

230	dependently interpretable factual and situational	279
231	units. Stage 2 identifies high-confidence implicit	280
232	information related to user preferences and profiles	281
233	without introducing new facts. Stage 3 performs	282
234	non-destructive normalization, including dedupli-	283
235	cation, coreference resolution, and temporal nor-	284
236	malization, producing aligned knowledge represen-	285
237	tations suitable for long-term storage. Each aligned	286
238	knowledge entry is stored as a note, represented as	287
239	a structured record containing an identifier, the ex-	288
240	tracted content, a semantic category, and associated	289
241	metadata.	290
242	<b>2.2.3 Knowledge Alignment</b>	291
243	To maintain semantic consistency across memory	292
244	layers, HiMem applies a unified alignment pro-	293
245	cess during memory construction. This process	294
246	includes temporal alignment of relative time ex-	295
247	pressions, coreference resolution for entity men-	296
248	tions, and extraction of implicit semantic relations.	297
249	Alignment is selectively applied: Episode Mem-	298
250	ory prioritizes preserving original dialogue context,	299
251	while Note Memory emphasizes abstraction and	300
252	normalization.	301
253	<b>2.3 Memory Retrieval</b>	302
254	HiMem supports two complementary retrieval	303
255	strategies. In <i>hybrid retrieval</i> , the system retrieves	304
256	information from both Episode Memory and Note	305
257	Memory to maximize recall. In contrast, <i>best-effort</i>	306
258	<i>retrieval</i> proceeds hierarchically by querying Note	307
259	Memory first and falling back to Episode Memory	308
260	only when evidence is deemed insufficient. Re-	309
261	trieved evidence is evaluated by an LLM to assess	310
262	answerability, and unsupported queries are explic-	
263	itly marked as unanswerable.	
264	<b>2.4 Memory Updating and Self-Evolution</b>	
265	During best-effort retrieval, HiMem employs a	
266	fixed LLM-based self-evaluation prompt to assess	
267	whether the retrieved evidence is sufficient to an-	
268	swer the query. This evaluation produces a binary	
269	judgment ( <i>sufficient</i> or <i>insufficient</i> ) under deter-	
270	ministic decoding (temperature = 0) and serves solely	
271	as a control signal, without introducing or revis-	
272	ing memory content. While related to iterative	
273	self-refinement approaches (Madaan et al., 2023),	
274	HiMem confines the LLM to deterministic routing	
275	and decision control.	
276	Memory reconsolidation is triggered only when	
277	two conditions are jointly satisfied: (i) retrieval	
278	from Note Memory alone is insufficient, and (ii)	
	the subsequently retrieved Episode Memory pro-	
	vides adequate supporting evidence. This conjunc-	
	tive trigger grounds updates in episodic context and	
	prevents premature revisions. Although conceptu-	
	ally related to reflective agent frameworks such as	
	(Shinn et al., 2023), HiMem performs structured,	
	evidence-grounded memory operations rather than	
	free-form verbal reflection.	
	When reconsolidation is triggered, HiMem con-	
	ducts query-conditioned knowledge extraction over	
	the supporting episodes and compares the extracted	
	information with existing notes. Their relationship	
	is classified as <i>independent</i> , <i>extendable</i> , or <i>contra-</i>	
	<i>dictory</i> , based on which the system applies ADD,	
	UPDATE, or DELETE operations to revise Note	
	Memory. This typed design avoids indiscriminate	
	overwriting and echoes classic belief revision per-	
	spectives (Gärdenfors, 1988), promoting long-term	
	stability and semantic consistency.	
	In contrast, episodic memory is treated as im-	
	mutable: newly constructed episodes are appended	
	chronologically without modification, preserving	
	the temporal integrity of interaction histories.	
	<b>2.5 Adaptive Forgetting</b>	
	To regulate memory growth under sustained in-	
	teractions, HiMem optionally employs an adap-	
	tive forgetting mechanism based on usage fre-	
	quency. In this work, forgetting primarily serves as	
	a scalability-oriented mechanism to control mem-	
	ory size and maintain retrieval efficiency, and does	
	not contribute directly to the performance gains	
	reported in our experiments.	
	<b>3 Experiments</b>	
	<b>3.1 Datasets</b>	
	We evaluate HiMem on <b>LoCoMo</b> (Maharana et al.,	
	2024), a benchmark designed to assess long-	
	horizon conversational reasoning under sustained	
	interactions. LoCoMo consists of multi-session di-	
	alogues with an average length of approximately	
	600 turns (around 16K tokens) and spans up to	
	32 interaction stages, posing significant challenges	
	for long-range dependency modeling and memory	
	management.	
	The benchmark covers diverse reasoning sce-	
	narios, including <i>Single-Hop</i> questions answerable	
	within a single session, <i>Multi-Hop</i> questions re-	
	quiring aggregation across distant dialogue turns,	
	<i>Temporal Reasoning</i> questions involving implicit	
	or explicit time relations, and <i>Open-Domain</i> ques-	

tions that combine dialogue content with external or commonsense knowledge. Following prior work (Chhikara et al., 2025), we exclude the Adversarial category from quantitative evaluation, as it focuses on unanswerability detection rather than answer correctness.

### 3.2 Evaluation Metrics

Since different long-term memory systems may apply varying degrees of compression or abstraction over dialogue histories, we adopt a multi-dimensional evaluation protocol to assess answer quality comprehensively. Specifically, following prior work that systematically studies LLM-as-a-Judge and its biases (Zheng et al., 2023; Pan et al., 2025; Xu et al., 2025b), we use GPT-4o-mini as the LLM judge to compute evaluation score (denoted as GPT-Score) as the primary metric to approximate semantic correctness and consistency, together with F1 to measure lexical overlap.

In addition, for efficiency evaluation, we report **latency** (Lat.) and **token consumption** (Tok.). Latency is measured as the time required for memory retrieval only, excluding LLM inference and response generation, in order to isolate the efficiency of the memory system.

### 3.3 Baselines

We compare HiMem with representative long-term memory frameworks that cover different design paradigms. **Mem0** (Chhikara et al., 2025) represents structured memory systems based on atomic factual extraction and graph-based organization. **SeCom** (Pan et al., 2025) adopts event-level semantic segmentation and compression to improve contextual completeness. **A-MEM** (Xu et al., 2025b) augments event-level memory with entities, relations, and temporal features to support time-aware retrieval and reasoning. These baselines enable a systematic comparison across retrieval-based, compressed-context, and structured memory approaches.

These baselines are selected based on their compatibility with long-horizon conversational memory, availability of reproducible implementations, and suitability for evaluation under a unified agent interface with comparable inference budgets; methods that primarily target system-level context management or non-dialogue memory access are therefore not included.

### 3.4 Settings

To ensure fair comparison, all methods are evaluated using the same base language model and identical decoding configurations. We use **GPT-4o-mini** as the backbone LLM and a shared embedding model for vector representations. For each evaluation setting, we conduct three independent trials with fixed prompts.

For baseline-comparative main results, we report mean $\pm$ std over multiple runs to reflect run-to-run variability. For auxiliary analyses and ablations (e.g., Table 3 and Table 4), we report mean values only for compact presentation, as these results are primarily intended to validate relative trends rather than to serve as headline comparisons.

Additional implementation details, including model configurations and hardware specifications, are provided in the Appendix.

## 4 Results and Analyses

### 4.1 Main Results

We evaluate HiMem on the LoCoMo benchmark to assess its ability to preserve, retrieve, and utilize information over long-horizon dialogues. Table 1 reports the performance of HiMem and representative baseline methods across diverse reasoning categories, including Single-Hop, Multi-Hop, Temporal Reasoning, and Open-Domain questions.

Overall, HiMem consistently outperforms all baselines across almost all categories. In particular, HiMem achieves substantial improvements on Multi-Hop and Temporal Reasoning tasks, which require aggregating scattered evidence across long interaction histories. These results indicate that hierarchical memory organization enables more effective modeling of long-range dependencies and semantic consistency than flat or monolithic memory structures. Moreover, the strong performance on Open-Domain questions suggests that HiMem can reliably integrate dialogue-derived knowledge with external or implicit information over extended time spans.

### 4.2 Ablation Study: Memory Components

To examine the contribution of different memory components, we conduct an ablation study by selectively removing Episode Memory or Note Memory from HiMem. The results are shown in Table 2.

Removing Episode Memory leads to a pronounced performance degradation across most categories, particularly on Multi-Hop and Temporal

Table 1: **Performance comparison of HiMem and baseline methods on LoCoMo.** Results are reported as mean (std) over three runs in percentages (%). Best results of GPT-Score are shown in **bold**, and second-best results are underlined.

Task	A-MEM		SeCom		Mem0		HiMem	
	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1
Single Hop	59.33 <sub>(0.51)</sub>	34.45 <sub>(0.46)</sub>	<u>87.02</u> <sub>(0.35)</sub>	23.70 <sub>(0.06)</sub>	75.90 <sub>(0.74)</sub>	53.05 <sub>(0.65)</sub>	<b>89.22</b> <sub>(0.06)</sub>	43.93 <sub>(0.24)</sub>
Multi Hop	40.78 <sub>(0.77)</sub>	20.98 <sub>(0.05)</sub>	<u>59.10</u> <sub>(1.17)</sub>	13.21 <sub>(0.01)</sub>	56.62 <sub>(2.86)</sub>	32.90 <sub>(1.11)</sub>	<b>70.92</b> <sub>(0.77)</sub>	28.32 <sub>(0.05)</sub>
Temporal	50.26 <sub>(1.55)</sub>	35.84 <sub>(0.26)</sub>	33.54 <sub>(0.39)</sub>	4.28 <sub>(0.06)</sub>	<u>68.54</u> <sub>(0.51)</sub>	56.37 <sub>(0.74)</sub>	<b>74.77</b> <sub>(0.25)</sub>	22.05 <sub>(0.22)</sub>
Open Domain	24.65 <sub>(2.14)</sub>	9.30 <sub>(0.50)</sub>	<b>60.07</b> <sub>(0.49)</sub>	8.57 <sub>(0.10)</sub>	42.36 <sub>(0.49)</sub>	22.70 <sub>(0.20)</sub>	<u>54.86</u> <sub>(1.30)</sub>	18.92 <sub>(0.45)</sub>
Overall	51.88 <sub>(0.52)</sub>	30.71 <sub>(0.29)</sub>	<u>69.03</u> <sub>(0.24)</sub>	16.77 <sub>(0.02)</sub>	68.74 <sub>(0.98)</sub>	48.16 <sub>(0.73)</sub>	<b>80.71</b> <sub>(0.21)</sub>	34.95 <sub>(0.11)</sub>

Table 2: **Ablation study on memory components in HiMem.** HiMem includes both *Note Memory* and *Episode Memory*; **w/o Note Memory** removes the Note Memory while retaining the Episode Memory; **w/o Episode Memory** removes the Episode Memory while retaining the Note Memory.

Task	HiMem		- w/o Episode		- w/o Note	
	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1
Single Hop	<b>89.22</b> <sub>(0.06)</sub>	43.93 <sub>(0.24)</sub>	76.50 <sub>(0.24)</sub>	41.09 <sub>(0.09)</sub>	<u>89.02</u> <sub>(0.20)</sub>	45.14 <sub>(0.03)</sub>
Multi Hop	<b>70.92</b> <sub>(0.77)</sub>	28.32 <sub>(0.05)</sub>	56.26 <sub>(0.33)</sub>	26.29 <sub>(0.12)</sub>	<u>70.33</u> <sub>(0.44)</sub>	26.13 <sub>(0.25)</sub>
Temporal	<b>74.77</b> <sub>(0.25)</sub>	22.05 <sub>(0.22)</sub>	68.12 <sub>(0.59)</sub>	23.65 <sub>(0.31)</sub>	<u>72.48</u> <sub>(0.39)</sub>	29.35 <sub>(0.16)</sub>
Open Domain	<b>54.86</b> <sub>(1.30)</sub>	18.92 <sub>(0.45)</sub>	48.26 <sub>(0.49)</sub>	22.58 <sub>(0.55)</sub>	<u>48.61</u> <sub>(0.98)</sub>	16.81 <sub>(0.15)</sub>
Overall	<b>80.71</b> <sub>(0.21)</sub>	34.95 <sub>(0.11)</sub>	69.29 <sub>(0.05)</sub>	33.59 <sub>(0.06)</sub>	<u>79.63</u> <sub>(0.22)</sub>	36.60 <sub>(0.02)</sub>

Reasoning tasks. This observation highlights the importance of preserving fine-grained contextual evidence aligned with the original interaction process. Without Episode Memory, the system struggles to recover detailed event-level information necessary for tracing complex reasoning chains over long dialogues.

In contrast, removing Note Memory results in a smaller but still consistent performance drop. This suggests that structured knowledge representations primarily serve to accelerate information localization and stabilize semantic anchors, while detailed contextual evidence remains indispensable for coverage and reasoning. Together, these findings demonstrate that Episode Memory and Note Memory play asymmetric yet complementary roles: effective long-term memory systems must balance information fidelity and abstraction rather than relying solely on either raw dialogue context or aggressively compressed representations.

### 4.3 Ablation Study: Knowledge Alignment

We further examine the role of the **Knowledge Alignment** module by comparing different memory types with and without a unified semantic alignment space. As shown in Table 3, disabling

Knowledge Alignment causes a pronounced performance drop for *Note Memory*, indicating that a unified semantic space is crucial for extraction-based memories that do not retain raw dialogue context. Such alignment substantially improves intent understanding and memory localization. In contrast, removing Knowledge Alignment for *Episode Memory* slightly improves performance, suggesting that when segmentation is well-structured, additional semantic fusion may dilute information inherent in raw dialogue. Overall, although extracted knowledge representations are more compact and explicit, they are more sensitive to implicit semantics and coreference resolution, and therefore benefit more from unified semantic alignment.

### 4.4 Memory Self-Evolution

During best-effort retrieval, when *Note Memory* fails to return self-validated results while *Episode Memory* provides sufficient evidence, HiMem triggers the **Memory Reconsolidation** mechanism. Specifically, the system performs query-conditioned information extraction over retrieved Episode Memory results and supplements missing knowledge in Note Memory through conflict detection and dynamic updating. As shown in Fig-

Table 3: **Ablation study of the Knowledge Alignment module.** For clarity, we denote Knowledge Alignment as KA. (1) **HiMem**: alignment applied only to Note Memory; (2) **HiMem w/o KA**: alignment disabled for both Note Memory and Episode Memory; (3) **Note Memory**: alignment applied to extracted knowledge; (4) **Note Memory w/o KA**: extracted knowledge retained without alignment; (5) **Episode Memory**: alignment applied to segmented episodes; (6) **Episode Memory w/o KA**: segmented episodes retained without alignment.

Method	Single Hop		Multi Hop		Temporal		Open Domain		Average	
	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1
HiMem	<b>89.22</b>	43.93	<b>70.92</b>	28.32	<b>74.77</b>	22.05	<b>54.86</b>	18.92	<b>80.71</b>	34.95
- w/o KA	87.51	43.75	69.86	28.53	75.18	28.14	52.08	15.96	79.50	35.98
Note Memory	<b>66.51</b>	35.88	<b>54.26</b>	24.92	<b>67.39</b>	19.89	<b>50.35</b>	23.88	<b>63.44</b>	29.79
- w/o KA	61.79	34.16	46.81	23.57	60.12	19.77	42.71	18.52	57.51	28.25
Episode Memory	88.31	45.53	65.25	24.76	71.55	29.08	<b>48.61</b>	18.09	78.12	36.59
- w/o KA	<b>89.02</b>	45.14	<b>70.33</b>	26.13	<b>72.48</b>	29.35	<b>48.61</b>	16.81	<b>79.63</b>	36.60

Table 4: **Ablation study of the Note Memory.** For clarity, we denote Knowledge Alignment as KA, and refer to enabling Memory Self-Evolution as ME.

Method	Single Hop		Multi Hop		Temporal		Open Domain		Average	
	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1
<b>Note Memory</b>										
w/o KA	61.79	34.16	46.81	23.57	60.12	19.77	42.71	18.52	57.51	28.25
+KA	66.51	35.88	54.26	24.92	67.39	19.89	50.35	23.88	63.44	29.79
+KA & +ME	<b>76.50</b>	<b>41.09</b>	<b>56.26</b>	<b>26.29</b>	<b>68.12</b>	<b>23.65</b>	<b>48.26</b>	<b>22.58</b>	<b>69.29</b>	<b>33.59</b>

ure 2 and Table 4, enabling Memory Self-Evolution improves Note Memory performance by approximately **5.85%**, which further leads to a slight overall performance gain of about **0.28%**. These results demonstrate that Memory Reconsolidation is an effective mechanism for enabling long-term memory self-evolution.

#### 4.5 Discussion of Extended Analyses

Additional analyses, including the retrieval strategies, hyperparameter sensitivity, and efficiency trade-offs, are provided in the Appendix.

## 5 Discussion

Long-horizon conversational agents require more than extended context windows or incremental memory accumulation; they critically depend on how information is structured, abstracted, and revised over time. The empirical results consistently demonstrate that hierarchical memory organization is a necessary condition for robust long-term reasoning rather than an optional architectural refinement. Episode Memory and Note Memory play asymmetric yet complementary roles: the former preserves fine-grained contextual evidence aligned

with the original interaction process, while the latter consolidates stable, high-frequency knowledge into compact representations that substantially reduce retrieval cost. The performance degradation observed when either component is removed confirms that effective long-term memory systems must balance fidelity and abstraction, instead of relying solely on raw dialogue context or aggressive compression.

Beyond static organization, our findings highlight that memory updating cannot be treated as a purely similarity-driven or append-only process. In long-term interactions, newly observed information often partially overlaps with, extends, or contradicts existing knowledge. HiMem’s conflict-aware Memory Reconsolidation explicitly distinguishes these cases and applies differentiated update strategies, which proves essential for maintaining semantic consistency over time. Importantly, the gains brought by memory self-evolution do not arise from heuristic rewriting, but from a conservative feedback loop between retrieval failure, episodic evidence inspection, and targeted knowledge supplementation.

Finally, the comparison between hybrid and best-

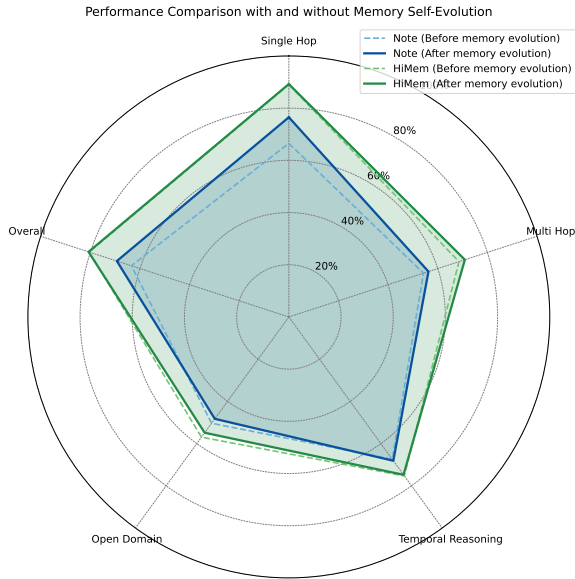


Figure 2: **Performance comparison (GPT-Score) before and after enabling Memory Self-Evolution for Note Memory and HiMem.** Memory Self-Evolution is triggered through conflict-aware *Memory Reconsolidation* during best-effort retrieval.

effort retrieval strategies indicates that hierarchical memory is not only a representational choice but also an efficiency mechanism. Retrieving abstract knowledge first and descending to concrete events only when necessary achieves a favorable trade-off between accuracy and computational cost, while simultaneously exposing latent information that can drive further memory evolution. Together, these observations suggest that long-horizon LLM agents should treat memory as a dynamic, multi-level system tightly coupled with retrieval and usage, rather than as a static external store.

## 6 Conclusion

This paper proposes **HiMem**, a hierarchical long-term memory framework for long-horizon dialogues, aimed at addressing several fundamental challenges faced by existing LLM agents in sustained interactions, including memory fragmentation, semantic drift, and the lack of self-evolution capability. Grounded in cognitive theories of human long-term memory, HiMem organically integrates event-level experiences with knowledge-level abstractions, and realizes efficient storage, retrieval, and dynamic updating of long-term information through a structured system design.

Methodologically, HiMem constructs cognitively consistent Episode Memory via **Topic-Aware Event-Surprise Dual-Channel Segmen-**

**tation**, providing fine-grained and semantically stable contextual support for complex reasoning tasks. Meanwhile, through a multi-stage information extraction pipeline and selective **Knowledge Alignment**, high-frequency and stable facts as well as user-specific attributes are consolidated into dense Note Memory representations, significantly reducing retrieval costs while preserving semantic fidelity. Furthermore, HiMem introduces a conflict-aware **Memory Reconsolidation** mechanism that closes the loop between retrieval and memory updating, enabling continuous correction and evolution of knowledge through usage.

Extensive experiments across multiple long-horizon conversational scenarios systematically validate the effectiveness of these design choices. HiMem consistently outperforms existing methods in terms of accuracy, temporal reasoning, and open-domain understanding. Ablation and analysis studies further reveal that these gains arise from the *synergistic interaction* among hierarchical memory structures, cognitively aligned event segmentation, memory-type-aware semantic alignment, and self-evolution mechanisms, rather than from isolated component-level improvements. In addition, analyses of retrieval modes and hyperparameters demonstrate that HiMem achieves a robust balance between knowledge coverage and system efficiency.

Overall, HiMem’s contributions extend beyond empirical performance improvements. More importantly, it offers a **practical paradigm for systematically integrating cognitive theories into the design of long-term memory for LLM agents**. By emphasizing memory-type distinctions, structured organization, and usage-driven feedback, this paradigm provides a methodological foundation for building scalable, interpretable, and self-evolving LLM agents. We hope that this work will inspire future research on long-term memory in more complex settings, including multi-agent, multimodal, and richly interactive environments.

## Limitations

Although HiMem demonstrates stable and significant performance advantages on long-horizon conversational tasks, several limitations remain that warrant further investigation. These limitations do not stem from flaws in the design itself, but rather reflect broader challenges commonly faced by long-term memory systems in realistic interac-

601 tive settings.

### 602 **Dependence on LLM Judgment Capabilities.**

603 First, HiMem relies extensively on the semantic  
604 and pragmatic judgment capabilities of the under-  
605 lying LLM during memory construction and up-  
606 dating, including event segmentation, information  
607 extraction, conflict detection, and evidence suffi-  
608 ciency evaluation. While experimental results indi-  
609 cate that such one-shot, rule-constrained judg-  
610 ments are stable and effective in practice, their  
611 quality inevitably depends on the capability of the  
612 base model. In scenarios involving noisy inputs,  
613 metaphorical language, or cross-cultural pragmatic  
614 variations, the accuracy of segmentation and knowl-  
615 edge extraction may be affected. Future work could  
616 explore incorporating lightweight auxiliary classi-  
617 fiers or uncertainty estimation mechanisms at criti-  
618 cal decision points to further enhance robustness  
619 under complex linguistic conditions.

### 620 **Expressive Limits of One-Shot Segmentation.**

621 Second, HiMem currently adopts a one-shot seg-  
622 mentation strategy, which offers clear advantages  
623 in efficiency and controllability, but also imposes  
624 an upper bound on expressive capacity. This strat-  
625 egy assumes that the event structure of a conversa-  
626 tion can be sufficiently identified through a single  
627 global pass. However, in extremely long or highly  
628 interleaved dialogues, event boundaries may ex-  
629 hibit hierarchical or recursive structures. Future  
630 extensions could investigate multi-granularity or it-  
631 erative event restructuring strategies, while preserv-  
632 ing the simplicity of the current design, to better  
633 accommodate non-linear conversational dynamics  
634 in Episode Memory.

### 635 **Conservative Triggers for Knowledge Evolution.**

636 Regarding memory self-evolution, HiMem primar-  
637 ily relies on retrieval failure or insufficient evi-  
638 dence as triggers for **Memory Reconsolidation**.  
639 While this conservative design promotes stability  
640 and avoids unnecessary updates, it may allow cer-  
641 tain latent inconsistencies or outdated knowledge  
642 to persist if they are not explicitly surfaced dur-  
643 ing retrieval. Designing more proactive yet noise-  
644 resistant evolution triggers remains an open chal-  
645 lenge. For example, future work could incorporate  
646 user feedback, cross-task consistency checks, or  
647 long-term statistical signals to detect and resolve  
648 implicit conflicts more effectively.

649 **Limited Evaluation Scope.** Finally, although  
650 HiMem is evaluated on representative long-horizon

651 dialogue benchmarks, the experiments are mainly  
652 confined to single-user, text-based interaction sce-  
653 narios. Real-world long-term interactions often in-  
654 volve multiple users, multimodal inputs, and richer  
655 social contexts, which impose additional demands  
656 on memory organization and updating. Extend-  
657 ing HiMem to multi-agent or multimodal settings,  
658 and studying how memories interact, conflict, and  
659 propagate across different agents, constitutes an  
660 important direction for future research.

661 Overall, these limitations highlight key research  
662 frontiers in advancing long-term memory systems  
663 from *usable* to truly *general-purpose*. HiMem pro-  
664 vides a viable pathway for systematically integrat-  
665 ing cognitive theories into the design of long-term  
666 memory for LLM agents. How to further enhance  
667 adaptability and generalization while maintaining  
668 structural clarity and interpretability remains a cen-  
669 tral focus for future work.

## 670 **Ethical Considerations**

671 We acknowledge that the development of hierar-  
672 chical long-term memory systems for LLM agents  
673 carries significant ethical responsibilities, particu-  
674 larly concerning data privacy, knowledge integrity,  
675 and potential societal impacts.

676 **Data Privacy and User Profiling** HiMem is de-  
677 signed to extract and store structured information,  
678 including user preferences and profiles, to maintain  
679 long-term interaction coherence. In real-world ap-  
680 plications, this involves the persistent storage of po-  
681 tentially sensitive personal information. We empha-  
682 size that any practical implementation of HiMem  
683 should adhere to privacy-by-design principles, such  
684 as the General Data Protection Regulation (GDPR).  
685 This includes implementing robust data encryption,  
686 ensuring transparency regarding what information  
687 is being "memorized," and providing users with the  
688 "right to be forgotten" by allowing them to inspect  
689 and delete specific entries in both Episode and Note  
690 Memory. In addition, the datasets used in this study  
691 are all publicly available and used in accordance  
692 with their respective licenses.

693 **Knowledge Integrity and Hallucinations** The  
694 "Memory Reconsolidation" mechanism introduces  
695 a dynamic self-evolution process where the system  
696 updates its internal knowledge based on new inter-  
697 actions. While this improves adaptability, it also  
698 poses a risk of "consolidating" hallucinations or  
699 incorrect information if the backbone LLM makes

erroneous judgments during the conflict-aware update phase. We have mitigated this through a conservative update strategy, but we caution that such systems should not be deployed in high-stakes domains (e.g., medical or legal advice) without human-in-the-loop verification.

**Bias Amplification** As HiMem relies on the semantic understanding and summarization capabilities of pre-trained LLMs, it may inadvertently inherit or amplify biases present in the foundation models during the memory abstraction process (Stage 1-3). We encourage future research to integrate bias-detection filters within the memory extraction pipeline to ensure that the "Notes" stored do not perpetuate harmful stereotypes or unfair social biases.

**Intended Use and Transparency** HiMem aims to foster more meaningful and efficient human-AI collaboration. However, the ability of an agent to form a "long-term bond" through persistent memory could potentially be misused for manipulative purposes. We advocate for full disclosure: users should be explicitly informed when they are interacting with an agent equipped with long-term memory capabilities to manage expectations and ensure informed consent.

## References

Nick Alonso, Tomás Figliolia, Anthony Ndirango, and Beren Millidge. 2024. [Toward conversational agents with context and time sensitive long-term memory](#). *Preprint*, arXiv:2406.00057.

Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet Singh, and Deshraj Yadav. 2025. [Mem0: Building production-ready AI agents with scalable long-term memory](#). *Preprint*, arXiv:2504.19413.

Yiming Du, Wenyu Huang, Danna Zheng, Zhaowei Wang, Sebastien Montella, Mirella Lapata, Kam-Fai Wong, and Jeff Z. Pan. 2025. [Rethinking memory in AI: Taxonomy, operations, topics, and future directions](#). *Preprint*, arXiv:2505.00675.

Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitan, Robert Osazuwa Ness, and Jonathan Larson. 2024. [From local to global: A graph RAG approach to query-focused summarization](#). *Preprint*, arXiv:2404.16130.

Gilles Fauconnier and Mark Turner. 2003. Conceptual blending, form and meaning. *Recherches en communication*, 19:57–86.

Zafeirios Fountas, Martin Benfeghoul, Adnan Omerjee, Fenia Christopoulou, Gerasimos Lampouras, Haitham Bou-Ammar, and Jun Wang. 2025. [Human-inspired episodic memory for infinite context LLMs](#). In *International Conference on Learning Representations (ICLR)*, Singapore. OpenReview.net.

Peter Gärdenfors. 1988. *Knowledge in Flux: Modeling the Dynamics of Epistemic States*. MIT Press.

Vanessa E. Ghosh and Asaf Gilboa. 2014. What is a memory schema? a historical perspective on current neuroscience literature. *Neuropsychologia*, 53:104–114.

Asaf Gilboa and Hannah Marlatte. 2017. Neurobiology of schemas and schema-mediated memory. *Trends in Cognitive Sciences*, 21(8):618–631.

Kostas Hatalis, Despina Christou, Joshua Myers, Steven Jones, Keith Lambert, Adam Amos-Binks, Zohreh Dannenhauer, and Dustin Dannenhauer. 2023. Memory matters: The need to improve long-term memory in LLM-agents. In *Proceedings of the AAAI Symposium Series*, volume 2, pages 277–280.

Zihong He, Weizhe Lin, Hao Zheng, Fan Zhang, Matt W. Jones, Laurence Aitchison, Xuhai Xu, Miao Liu, Per Ola Kristensson, and Junxiao Shen. 2024. [Human-inspired perspectives: A survey on AI long-term memory](#). *Preprint*, arXiv:2411.00489.

Xun Jiang, Feng Li, Han Zhao, Jiahao Qiu, Jiaying Wang, Jun Shao, Shihao Xu, Shu Zhang, Weiling Chen, Xavier Tang, Yize Chen, Mengyue Wu, Weizhi Ma, Mengdi Wang, and Tianqiao Chen. 2024. [Long-term memory: The foundation of AI self-evolution](#). *Preprint*, arXiv:2410.15665.

Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2024. [HippoRAG: Neurobiologically inspired long-term memory for large language models](#). In *Advances in Neural Information Processing Systems*.

Bernal Jiménez Gutiérrez, Yiheng Shu, Weijian Qi, Sizhe Zhou, and Yu Su. 2025. [From RAG to memory: Non-parametric continual learning for large language models](#). In *International Conference on Machine Learning (ICML)*, Vancouver, Canada. OpenReview.net.

Kuang-Huei Lee, Xinyun Chen, Hiroki Furuta, John F. Canny, and Ian Fischer. 2024. [A human-inspired reading agent with gist memory of very long contexts](#). In *International Conference on Machine Learning (ICML)*, Vienna, Austria. OpenReview.net.

Yubo Li, Xiaobin Shen, Xinyu Yao, Xueying Ding, Yidi Miao, Ramayya Krishnan, and Rema Padman. 2025. [Beyond single-turn: A survey on multi-turn interactions with large language models](#). *Preprint*, arXiv:2504.04717.

802	Jiahao Liu, Shengkang Gu, Dongsheng Li, Guangping	Hongjin Qian, Zheng Liu, Peitian Zhang, Kelong Mao,	859
803	Zhang, Mingzhe Han, Hansu Gu, Peng Zhang, Tun	Defu Lian, Zhicheng Dou, and Tiejun Huang. 2025.	860
804	Lu, Li Shang, and Ning Gu. 2025. <a href="#">AgentCF++:</a>	<a href="#">MemoRAG: Boosting long context processing with</a>	861
805	<a href="#">Memory-enhanced LLM-based agents for popularity-</a>	<a href="#">global memory-enhanced retrieval augmentation.</a>	862
806	<a href="#">aware cross-domain recommendations.</a>	<i>In Proceedings of The Web Conference (WWW 2025),</i>	863
807	<i>In Proceedings of the 48th International ACM SIGIR Confer-</i>	<i>pages 2366–2377, Sydney, Australia. ACM.</i>	864
808	<i>ence on Research and Development in Information</i>		
809	<i>Retrieval (SIGIR 2025), pages 2566–2571, Padua,</i>		
810	<i>Italy. ACM.</i>		
811	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	Alireza Rezaadeh, Zichao Li, Wei Wei, and Yujia Bao.	865
812	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	2025. <a href="#">From isolated conversations to hierarchical</a>	866
813	Nouha Dziri, Shrimai Prabhunoye, Yiming Yang,	<a href="#">schemas: Dynamic tree memory representation for</a>	867
814	Shashank Gupta, Bodhisattwa Prasad Majumder,	<a href="#">LLMs.</a>	868
815	Katherine Hermann, Sean Welleck, Amir Yazdan-	<i>In International Conference on Learning</i>	869
816	bakhsh, and Peter Clark. 2023. <a href="#">Self-refine: Iterative</a>	<i>Representations (ICLR), Singapore. OpenReview.net.</i>	
817	<a href="#">refinement with self-feedback.</a>		
818	<i>In Advances in Neural Information Processing Systems.</i>		
819	Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov,	Rana Salama, Jason Cai, Michelle Yuan, Anna Currey,	870
820	Mohit Bansal, Francesco Barbieri, and Yuwei Fang.	Monica Sunkara, Yi Zhang, and Yassine Benajiba.	871
821	2024. <a href="#">Evaluating very long-term conversational</a>	2025. <a href="#">MemInsight: Autonomous memory augmen-</a>	872
822	<a href="#">memory of LLM agents.</a>	<a href="#">tation for LLM agents.</a>	873
823	<i>In Proceedings of the 62nd Annual Meeting of the Association for Computational</i>	<i>Preprint, arXiv:2503.21760.</i>	
824	<i>Linguistics (Volume 1: Long Papers), pages 13851–</i>		
825	<i>13870, Bangkok, Thailand. Association for Computa-</i>		
826	<i>tational Linguistics.</i>		
827	Karim Nader. 2015. Reconsolidation and the dynamic	Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh	874
828	nature of memory. <i>Cold Spring Harbor Perspectives</i>	Khanna, Anna Goldie, and Christopher D. Manning.	875
829	<i>in Biology</i> , 7(10):a021782.	2024. <a href="#">RAPTOR: Recursive abstractive processing</a>	876
830		<a href="#">for tree-organized retrieval.</a>	877
831		<i>In International Conference on Learning Representations (ICLR), Vienna,</i>	878
832		<i>Austria. OpenReview.net.</i>	879
833			
834		Sentence-Transformers. 2021. <a href="#">sentence-</a>	880
835		<a href="#">transformers/all-mpnet-base-v2.</a>	881
836		Hugging Face	882
837		model card.	
838			
839		Lianlei Shan, Shixian Luo, Zezhou Zhu, Yu Yuan, and	883
840		Yong Wu. 2025. <a href="#">Cognitive memory in large language</a>	884
841		<a href="#">models.</a>	885
842		<i>Preprint, arXiv:2504.02441.</i>	
843			
844		Noah Shinn, Federico Cassano, Ashwin Gopinath,	886
845		Karthik Narasimhan, and Shunyu Yao. 2023. <a href="#">Re-</a>	887
846		<a href="#">flexion: Language agents with verbal reinforcement</a>	888
847		<a href="#">learning.</a>	889
848		<i>In Advances in Neural Information Process-</i>	890
849		<i>ing Systems.</i>	
850			
851		Xiaorui Su, Yibo Wang, Shanghua Gao, Xiaolong	891
852		Liu, Valentina Giunchiglia, Djork-Arné Clevert, and	892
853		Marinka Zitnik. 2025. <a href="#">KGAREvion: An AI agent</a>	893
854		<a href="#">for knowledge-intensive biomedical QA.</a>	894
855		<i>In International Conference on Learning Representations</i>	895
856		<i>(ICLR), Singapore. OpenReview.net.</i>	896
857			
858		Zhen Tan, Jun Yan, I-Hung Hsu, Rujun Han, Zifeng	897
		Wang, Long T. Le, Yiwen Song, Yanfei Chen, Hamid	898
		Palangi, George Lee, Anand Rajan Iyer, Tianlong	899
		Chen, Huan Liu, Chen-Yu Lee, and Tomas Pfister.	900
		2025. <a href="#">In prospect and retrospect: Reflective mem-</a>	901
		<a href="#">ory management for long-term personalized dialogue</a>	902
		<a href="#">agents.</a>	903
		<i>In Proceedings of the 63rd Annual Meet-</i>	904
		<i>ing of the Association for Computational Linguistics</i>	905
		<i>(Volume 1: Long Papers), pages 8416–8439, Vienna,</i>	906
		<i>Austria. Association for Computational Linguistics.</i>	
		Yaxiong Wu, Sheng Liang, Chen Zhang, Yichao Wang,	907
		Yongyue Zhang, Huifeng Guo, Ruiming Tang, and	908
		Yong Liu. 2025a. <a href="#">From human memory to AI mem-</a>	909
		<a href="#">ory: A survey on memory mechanisms in the era of</a>	910
		<a href="#">LLMs.</a>	911
		<i>Preprint, arXiv:2504.15965.</i>	
		Yaxiong Wu, Yongyue Zhang, Sheng Liang, and	912
		Yong Liu. 2025b. <a href="#">SGMem: Sentence graph mem-</a>	913
		<a href="#">ory for long-term conversational agents.</a>	914
		<i>Preprint,</i>	915
		<i>arXiv:2509.21212.</i>	

916 Derong Xu, Yi Wen, Pengyue Jia, Yingyi Zhang, Wen-  
917 lin Zhang, Yichao Wang, Huifeng Guo, Ruiming  
918 Tang, Xiangyu Zhao, Enhong Chen, and Tong Xu.  
919 2025a. [Towards multi-granularity memory associa-](#)  
920 [tion and selection for long-term conversational agents.](#)  
921 *Preprint*, arXiv:2505.19549.

922 Jing Xu, Arthur Szlam, and Jason Weston. 2022. [Be-](#)  
923 [yond goldfish memory: Long-term open-domain con-](#)  
924 [versation.](#) In *Proceedings of the 60th Annual Meet-*  
925 [ing of the Association for Computational Linguistics](#)  
926 [\(Volume 1: Long Papers\)](#), pages 5180–5197, Dublin,  
927 Ireland. Association for Computational Linguistics.

928 Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao,  
929 Juntao Tan, and Yongfeng Zhang. 2025b. [A-](#)  
930 [MEM: Agentic memory for LLM agents.](#) *Preprint*,  
931 arXiv:2502.12110.

932 Yangyang Yu, Haohang Li, Zhi Chen, Yuechen Jiang,  
933 Yang Li, Jordan W. Suchow, Denghui Zhang,  
934 and Khaldoun Khashanah. 2025. [FinMem: A](#)  
935 [performance-enhanced LLM trading agent with lay-](#)  
936 [ered memory and character design.](#) *IEEE Transac-*  
937 *tions on Big Data*, 11(6):3443–3459.

938 Xihang Yue, Linchao Zhu, and Yi Yang. 2024. [FragRel:](#)  
939 [Exploiting fragment-level relations in the external](#)  
940 [memory of large language models.](#) In *Findings of*  
941 [the Association for Computational Linguistics: ACL](#)  
942 [2024](#), pages 16348–16361, Bangkok, Thailand. As-  
943 sociation for Computational Linguistics.

944 Guibin Zhang, Muxin Fu, Guancheng Wan, Miao Yu,  
945 Kun Wang, and Shuicheng Yan. 2025a. [G-Memory:](#)  
946 [Tracing hierarchical memory for multi-agent systems.](#)  
947 *Preprint*, arXiv:2506.07398.

948 Yujie Zhang, Weikang Yuan, and Zhuoren Jiang. 2025b. [Bridging](#)  
949 [intuitive associations and deliberate recall:](#)  
950 [Empowering LLM personal assistant with graph-](#)  
951 [structured long-term memory.](#) In *Findings of the As-*  
952 [sociation for Computational Linguistics: ACL 2025](#),  
953 pages 17533–17547, Vienna, Austria. Association  
954 for Computational Linguistics.

955 Zeyu Zhang, Quanyu Dai, Xiaohu Bo, Chen Ma, Rui Li,  
956 Xu Chen, Jieming Zhu, Zhenhua Dong, and Ji-Rong  
957 Wen. 2025c. [A survey on the memory mechanism of](#)  
958 [large language model-based agents.](#) *ACM Transac-*  
959 *tions on Information Systems*, 43(6):155:1–155:47.

960 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan  
961 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,  
962 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,  
963 Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging](#)  
964 [LLM-as-a-judge with MT-bench and chatbot arena.](#)  
965 In *Advances in Neural Information Processing Sys-*  
966 *tems*.

967 Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and  
968 Yanlin Wang. 2024. [MemoryBank: Enhancing large](#)  
969 [language models with long-term memory.](#) In *Pro-*  
970 [ceedings of the AAAI Conference on Artificial Intelli-](#)  
971 [gence](#), volume 38, pages 19724–19731.

## A Related Work 972

973 In recent years, research on long-term memory  
974 for LLM agents has progressed along multiple  
975 technical directions in parallel, including retrieval-  
976 augmented generation, long-context modeling, and  
977 structured memory systems. These paradigms re-  
978 spectively emphasize external information access,  
979 context capacity expansion, and long-term knowl-  
980 edge organization. However, most existing ap-  
981 proaches are developed from isolated design per-  
982 spectives and lack a unified abstraction to sys-  
983 tematically characterize the commonalities, dif-  
984 ferences, and inherent trade-offs among different  
985 memory mechanisms. To this end, we introduce  
986 a three-dimensional analytical framework, termed  
987 **Memory Form–Memory Organization–Memory**  
988 **Operation**, which revisits the design space of  
989 LLM long-term memory systems from the perspec-  
990 tives of memory unit representation, organizational  
991 structure, and dynamic operations, providing a uni-  
992 fied basis for comparing different research direc-  
993 tions.

### A.1 Memory Form 994

995 **Memory Form** describes the fundamental repre-  
996 sentation and granularity of memory units, deter-  
997 mining their content structure and serving as the  
998 foundational component of long-term memory sys-  
999 tems. Early approaches predominantly rely on  
1000 static segmentation strategies based on dialogue  
1001 turns or sessions (Tan et al., 2025; Wu et al., 2025b;  
1002 Salama et al., 2025), which often fail to align with  
1003 semantic boundaries or the temporal evolution of  
1004 events. Recent work such as SeCom (Pan et al.,  
1005 2025) introduces semantic and event boundary de-  
1006 tection mechanisms (Event Segmentation) to seg-  
1007 ment dialogues at the semantic level, leading to no-  
1008 table improvements in semantic coherence and con-  
1009 textual completeness. MemGAS (Xu et al., 2025a)  
1010 further provides multiple segmentation modes, al-  
1011 lowing systems to select different memory granu-  
1012 larities according to task objectives.

1013 Regarding memory content construction, SeCom  
1014 builds structured memory representations through  
1015 event segmentation and semantic summarization,  
1016 while Mem0 constructs fragmented fact units via  
1017 information extraction. A-MEM (Xu et al., 2025b)  
1018 and THEANINE (Ong et al., 2025) further enrich  
1019 memory representations by incorporating multi-  
1020 dimensional features such as entities, relations, and  
1021 temporal attributes. Collectively, these approaches

1022 explore the balance between information complete- 1073  
1023 ness and noise suppression, while supporting tem- 1074  
1024 poral reasoning through timestamps or explicit tem- 1075  
1025 poral modeling. For instance, A-MEM preserves 1076  
1026 timestamps to track event order, whereas THEA- 1077  
1027 NINE explicitly models temporal dependencies 1078  
1028 among memory units to capture dynamic seman- 1079  
1029 tics. 1080

1030 Despite these advances, existing methods largely 1081  
1031 focus on explicit semantic segmentation and static 1082  
1032 feature encoding, without sufficiently modeling im- 1083  
1033 plicit semantic dependencies or hierarchical inter- 1084  
1034 actions among memory units. In contrast, long- 1085  
1035 context modeling approaches such as MemGPT 1086  
1036 (Packer et al., 2023) approach the problem from a 1087  
1037 system capacity management perspective, empha- 1088  
1038 sizing contextual continuity and memory schedul- 1089  
1039 ing rather than semantic structuring of memory 1090  
1040 content. As a result, they exhibit limitations in fine- 1091  
1041 grained factual recall and long-term consistency 1092  
1042 modeling. 1093

1043 Unlike prior work, HiMem refines event gran- 1094  
1044 ularity through **Dual-Channel Segmentation** 1095  
1045 and integrates semantic-level fusion with multi- 1096  
1046 dimensional feature encoding, unifying explicit 1097  
1047 structural representations with implicit semantic 1098  
1048 modeling. This design reflects a shift in mem- 1099  
1049 ory form from structure-centric storage toward 1100  
1050 semantic-alignment-centric modeling.

## 1051 A.2 Memory Organization 1100

1052 **Memory Organization** characterizes how mem- 1101  
1053 ory units are connected and organized, directly af- 1102  
1054 fecting retrieval efficiency and scalability during 1103  
1055 reasoning. Existing approaches generally follow 1104  
1056 two main directions. On one hand, structured mem- 1105  
1057 ory organizations are adopted by methods such 1106  
1058 as SeCom and A-MEM, which connect memory 1107  
1059 units linearly along temporal or topical dimen- 1108  
1060 sions to maintain semantic continuity. Some frame- 1109  
1061 works (Rezazadeh et al., 2025; Zhang et al., 2025b; 1110  
1062 Chhikara et al., 2025; Zhang et al., 2025a) fur- 1111  
1063 ther introduce tree or graph structures to capture 1112  
1064 cross-event and cross-topic semantic relations. On 1113  
1065 the other hand, some approaches draw inspiration 1114  
1066 from human cognition or operating systems (He 1115  
1067 et al., 2024; Shan et al., 2025). For example, Hip- 1116  
1068 poRAG2 (Jiménez Gutiérrez et al., 2025) combines 1117  
1069 graph structures with vector spaces to simulate hip- 1118  
1070 pocampal indexing mechanisms, enhancing seman- 1119  
1071 tic association and retrieval accuracy. MemGPT 1120  
1072 mimics page caching mechanisms in operating sys-

1073 tems, organizing memory hierarchically based on 1074  
1075 access cost and treating the context window as a 1076  
1077 high-speed working memory. 1078

1079 However, these methods typically seek com- 1080  
1081 promises within a single organizational structure, 1081  
1082 making it difficult to simultaneously accommodate 1082  
1083 diverse task requirements and storage efficiency. 1083  
1084 Moreover, they primarily focus on data structure 1084  
1085 design while paying limited attention to explicitly 1085  
1086 modeling hierarchical differences in memory con- 1086  
1087 tent. To address this limitation, HiMem adopts a 1087  
1088 multi-level hybrid organization strategy that con- 1088  
1089 structs a hierarchical long-term memory structure 1089  
1090 spanning from concrete events to abstract knowl- 1090  
1091 edge, based on the information density and abstrac- 1091  
1092 tion level of memory units. During retrieval, seman- 1092  
1093 tic filtering progressively narrows down relevant 1093  
1094 information, significantly reducing computational 1094  
1095 cost and noise while preserving high precision. 1095

1096 Beyond structural optimization, HiMem’s key 1096  
1097 contribution lies in achieving **Hierarchical Decou- 1097  
1098 pling** in memory organization. By jointly opti- 1098  
1099 mizing semantic association, knowledge indexing, 1099  
1100 and retrieval efficiency within a unified hierarchi- 1100  
1101 cal framework, HiMem transitions from structure- 1101  
1102 driven organization to semantic-driven organiza- 1102  
1103 tion. 1103

## 1104 A.3 Memory Operation 1100

1105 **Memory Operation** focuses on dynamic memory 1101  
1106 updates and operational mechanisms during usage, 1102  
1107 which are critical for long-term adaptability and 1103  
1108 self-evolution. Self-evolving memory is widely re- 1104  
1109 garded as a core capability for long-horizon LLM 1105  
1110 agents, enabling systems to continuously learn, up- 1106  
1111 date, and refine their knowledge structures through 1107  
1112 sustained interactions (Jiang et al., 2024). Sev- 1108  
1113 eral works, including Mem0, A-MEM, and THEA- 1109  
1114 NINE, support dynamic updates of memory units to 1110  
1115 reflect new conversational content. MemoryBank 1111  
1116 (Zhong et al., 2024) introduces a forgetting-curve- 1112  
1117 based decay mechanism that periodically removes 1113  
1118 low-frequency or irrelevant information to reduce 1114  
1119 storage pressure and semantic interference. 1115

1120 MemGPT emphasizes self-management oper- 1116  
1121 ations of memory: the LLM dynamically reads, 1117  
1122 writes, and schedules memory across different stor- 1118  
1123 age layers based on task requirements, triggering 1119  
1124 paging and summarization when the context win- 1120  
1125 dow is constrained, and writing historical infor- 1121  
1126 mation to external storage to maintain contextual 1122  
1127 continuity. Notably, unlike systems that achieve 1123

Method	Memory Form	Degeneracy	Memory Organization	Degeneracy	Memory Operation	Degeneracy
SeCom	Event-level semantic summaries	Yes	Linear/hybrid index or organized by time	Yes	Append-only memory construction	Yes
A-MEM	Zettelkasten-style notes (note-level abstraction)	Partial	Linked note graph with local semantic connections	Partial	Incremental update and merge without explicit conflict typing	Yes
Mem0	Atomic factual triples	Yes	Graph-based memory with semantic edges	Partial	Update and deletion based on similarity and recency	Yes
MemGPT	Page-based memory blocks constrained by context window	Yes	External memory swapping driven by capacity management	Yes	Eviction and replacement without semantic consistency modeling	Yes
HiMem	<b>Hierarchical Episode → Knowledge</b>	<b>No</b>	<b>Semantic-driven hierarchical decoupling</b>	<b>No</b>	<b>Conflict-aware reconsolidation (assimilation/accommodation)</b>	<b>No</b>

Figure 3: **Mapping of representative long-term memory systems under the Memory Form–Memory Organization–Memory Operation framework.** The framework characterizes long-term memory systems along three dimensions: memory unit representation, organizational structure, and memory operations. When a dimension collapses into a single fixed design choice that restricts adaptive trade-offs among granularity, structure, or temporal evolution, it is considered to exhibit design degeneration. In contrast, HiMem maintains non-degenerate designs across all three dimensions, enabling more flexible and evolvable long-term memory modeling.

self-evolution through content-level updates or reconstructions, MemGPT does not directly modify the semantic structure of memory, but instead manages limited context resources via operating-system-style scheduling and compression.

Nevertheless, although these methods support dynamic updates to some extent, they largely remain at the level of content addition or deletion, lacking mechanisms for conflict awareness and semantic reintegration at the knowledge level. When semantic conflicts arise between new and existing information, systems often struggle to decide whether to retain, merge, or revise memories, leading to degraded knowledge consistency or increased computational overhead due to excessive filtering and rewriting.

To address this limitation, HiMem introduces a conflict-aware dynamic evolution mechanism based on memory type differentiation. For memory units that record objective events, whose semantics are relatively stable, conflict detection is unnecessary. In contrast, for knowledge-oriented memories representing user preferences or personal traits, semantic conflicts trigger assimilation or accommodation operations, enabling Memory Reconsolidation and supporting self-correction and continuous evolution. This design closely aligns with cognitive theories of Memory Reconsolidation (Nader, 2015), allowing LLM agents to maintain consistency during long-term interactions.

As illustrated in Figure 3, we map representative long-term memory systems into the proposed three-dimensional analytical framework. Most existing approaches avoid design degeneration in only

one dimension, while implicitly simplifying the remaining dimensions into fixed design choices. For instance, some methods improve memory form through event segmentation or summarization, but retain static assumptions in memory organization and updating mechanisms. Others enhance organizational flexibility or operational strategies, yet remain constrained in memory form and semantic modeling.

In contrast, HiMem is among the few systems that maintain non-degenerate designs across all three dimensions. It introduces a hierarchical “event-to-knowledge” memory form, adopts semantic hierarchy rather than static topology as the organizing principle, and explicitly models conflict-aware reconsolidation in memory operations. This three-dimensional non-degeneracy enables HiMem to support more adaptive and evolvable memory management in long-term interaction scenarios.

## B Implementation Details

To ensure fair and stable comparisons, all methods are evaluated under identical experimental settings. We use **GPT-4o-mini** as the base language model for all methods with fixed decoding parameters (temperature=0.0, max\_tokens=8192), set top-k=10 for memory selection during retrieval, and adopt **all-mpnet-base-v2** (Sentence Transformers, 2021) for vector representations. All remaining implementation details strictly follow the official implementations and recommended configurations of the corresponding baselines.

For each benchmark evaluation, we keep model parameters and prompts fixed, run three repeated

evaluations, and report the mean and standard deviation of all key metrics across runs. We apply the same evaluation protocol to HiMem and all baselines. In addition, we use GPT-4o-mini as the judge for GPT-Score.

All experiments are conducted on the same hardware environment—a MacBook Pro with Apple M4 Max and 128GB unified memory—to eliminate hardware differences as a confounding factor for efficiency metrics.

For HiMem, we use **OpenSearch** (OpenSearch Project, 2025) as the storage backend for Episode Memory, and adopt **Qdrant** (Qdrant Team, 2025) to manage knowledge-oriented memory representations in Note Memory.

## C Extended Analyses

### C.1 Memory Retrieval Modes

We compare two retrieval strategies in HiMem: **hybrid retrieval**, which simultaneously queries Note Memory and Episode Memory, and **best-effort retrieval**, which first queries Note Memory and descends to Episode Memory only when necessary. As shown in Table 5, hybrid retrieval achieves high accuracy while significantly reducing response latency, indicating that the hierarchical memory structure enables parallel localization of relevant information and thus shortens reasoning time. However, as shown in Table 6, hybrid retrieval incurs the highest token consumption due to aggregating information from both memory layers. In contrast, best-effort retrieval introduces additional latency due to self-evaluation and potential lower-layer queries, but consumes substantially fewer tokens since most queries are resolved at the Note Memory level. These results validate that hierarchical retrieval from abstract knowledge to concrete events achieves an effective trade-off between knowledge coverage and system efficiency.

### C.2 Hyperparameter Analysis

We further analyze the effect of the top-k parameter ( $k \in \{5, 10, 15, 20, 25\}$ ) on system performance. The results show that increasing  $k$  improves retrieval coverage and accuracy, but performance plateaus when  $k \geq 10$ . Meanwhile, retrieval latency and token consumption increase with larger  $k$ , reflecting higher retrieval costs. These findings indicate that HiMem, through fine-grained Episode segmentation and multi-stage knowledge extraction, can capture sufficient information within a

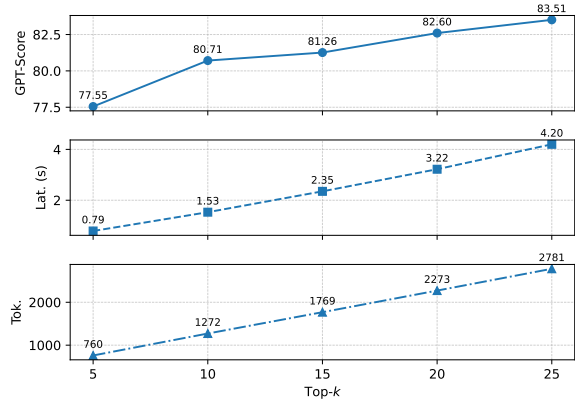


Figure 4: **Effect of top-k on performance and efficiency.** We report GPT-Score, latency, and token consumption as a function of  $k$ . Latency reports search latency only.

small retrieval window, achieving optimal performance without expanding the search scope. In contrast, excessively large  $k$  introduces irrelevant information and unnecessary processing overhead. This analysis further confirms the efficiency and information-density advantages of HiMem in long-horizon retrieval.

### C.3 Efficiency Analysis

We evaluate efficiency by comparing response latency and token consumption across methods. Although HiMem introduces semantic fusion and hierarchical retrieval, its latency is not significantly worse than that of other baselines, indicating that the computational overhead of lightweight semantic alignment is minimal. In contrast, by leveraging accurate intent modeling and prior knowledge retrieval, HiMem can more rapidly localize required information. In terms of token consumption, HiMem substantially outperforms all baselines, benefiting from its hierarchical memory structure that prioritizes highly condensed knowledge units and thereby reduces context length and generation burden.

## D Discussion

The experimental results systematically demonstrate that HiMem’s advantages in long-horizon conversational tasks do not stem from isolated improvements in individual components. Instead, they arise from a set of *interdependent design decisions* spanning memory representation, organization, and evolution. Together, these decisions form a closed loop from *information acquisition* to

Table 5: Performance comparison of HiMem under hybrid and best-effort retrieval strategies.

Strategy	Single Hop		Multi Hop		Temporal		Open Domain		Average	
	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1	GPT-Score	F1
hybrid retrieval	<b>89.22</b>	43.93	<b>70.92</b>	28.32	<b>74.77</b>	22.05	<b>54.86</b>	18.92	<b>80.71</b>	34.95
best-effort retrieval	83.59	43.93	62.88	27.42	72.38	25.26	47.92	20.34	75.24	35.54

Table 6: Efficiency comparison of HiMem under hybrid and best-effort retrieval strategies, measured by latency and token consumption.

Strategy	Single Hop		Multi Hop		Temporal		Open Domain		Average	
	Lat.(s)	Tok.	Lat.(s)	Tok.	Lat.(s)	Tok.	Lat.(s)	Tok.	Lat.(s)	Tok.
hybrid retrieval	1.57	1292.53	1.55	1292.17	1.35	1177.69	1.67	1343.25	1.53	1271.69
best-effort retrieval	1.63	1016.45	2.03	1257.29	1.88	1200.50	2.66	1583.06	1.82	1134.24

1272 *knowledge consolidation* and further to *continuous*  
 1273 *correction*, enabling stable performance in complex  
 1274 and dynamic conversational environments.

### 1275 **D.1 Hierarchical Memory as a Fundamental** 1276 **Constraint for Long-Term Dialogue** 1277 **Modeling**

1278 Both the main results and ablation studies con-  
 1279 sistentlly indicate that the **hierarchical memory**  
 1280 **structure** (Episode Memory + Note Memory) con-  
 1281 stitutes the core prerequisite for long-term conver-  
 1282 sational performance. Episode Memory preserves  
 1283 fine-grained contextual information aligned with  
 1284 the original interaction process, allowing the sys-  
 1285 tem to accurately trace back critical evidence in  
 1286 tasks such as multi-hop reasoning and temporal  
 1287 dependency modeling. In contrast, Note Memory  
 1288 compresses high-frequency and stable knowledge  
 1289 into dense semantic units through information ex-  
 1290 traction and structured representation, substantially  
 1291 reducing retrieval cost.

1292 The functional asymmetry between these two  
 1293 memory types explains the patterns observed in  
 1294 the ablation results. Removing Episode Memory  
 1295 leads to a severe performance drop, highlighting  
 1296 that *raw contextual information remains indispens-*  
 1297 *able* for complex reasoning tasks. Removing Note  
 1298 Memory also degrades performance, but to a lesser  
 1299 extent, indicating that structured knowledge pri-  
 1300 marily serves as a mechanism for accelerating lo-  
 1301 calization and stabilizing semantic anchors. These  
 1302 findings underscore a key insight: **the effectiveness**  
 1303 **of long-term memory systems lies not in com-**  
 1304 **plete abstraction, but in maintaining a balance**

**between abstraction and fidelity.** 1305

### 1306 **D.2 The Decisive Role of Dual-Channel** 1307 **Segmentation in Modeling Long-Term** 1308 **Dependencies**

1309 The effectiveness of Episode Memory further de-  
 1310 pends on how it is constructed. Experimental  
 1311 results show that Episode units built via **Topic-**  
 1312 **Aware Event-Surprise Dual-Channel Segmen-**  
 1313 **tation** consistently outperform coarse-grained or  
 1314 purely topic-based segmentation strategies in long-  
 1315 term retrieval and reasoning. This suggests that  
 1316 semantic continuity alone is insufficient to capture  
 1317 event boundaries in real conversations; *cognitive-*  
 1318 *level discontinuities*, such as shifts in emotion, in-  
 1319 tent, or discourse function, are equally critical sig-  
 1320 nals for long-term memory modeling.

1321 By explicitly incorporating both “topic conti-  
 1322 nuity” and “surprise-driven discontinuity” crite-  
 1323 ria at the segmentation stage and fusing them via  
 1324 an OR rule, HiMem generates memory units that  
 1325 better align with human event perception. Such  
 1326 cognitively consistent segmentation reduces cross-  
 1327 segment interference and improves the likelihood  
 1328 that retrieval targets genuinely relevant contexts,  
 1329 which directly manifests as performance gains in  
 1330 Multi-Hop and Temporal Reasoning tasks.

### 1331 **D.3 Selective Effects of Knowledge Alignment** 1332 **Reveal Memory-Type Differences**

1333 The ablation study on Knowledge Alignment  
 1334 uncovers a more nuanced and informative phe-  
 1335 nomenon: **a unified semantic alignment space**  
 1336 **does not benefit all memory types equally.** For  
 1337 Note Memory, removing Knowledge Alignment re-

Table 7: **Efficiency comparison of HiMem and baseline methods on the LoCoMo dataset.** The lowest latency and smallest token consumption are highlighted in bold, while the second-best results are underlined. Latency (Lat.) reports search latency only and is measured in seconds (s). Notably, SeCom achieves lower latency by preloading per-sample data into memory prior to inference; therefore, its latency is not directly comparable to that of other methods.

Task	A-MEM		SeCom		Mem0		HiMem	
	Lat.(s)	Tok.	Lat.(s)	Tok.	Lat.(s)	Tok.	Lat.(s)	Tok.
Single Hop	<b>0.93</b> (0.08)	2698.35(1.20)	-	2738.62(0.05)	4.65(0.13)	<u>1586.37</u> (207.51)	<u>1.57</u> (0.03)	<b>1292.53</b> (0.03)
Multi Hop	<b>0.91</b> (0.14)	2715.48(3.35)	-	2742.68(0.19)	3.96(0.19)	<u>1588.96</u> (209.64)	<u>1.55</u> (0.03)	<b>1292.17</b> (0.15)
Temporal	<b>0.95</b> (0.12)	2697.38(1.43)	-	2612.02(0.05)	4.64(0.22)	<u>1591.95</u> (209.63)	<u>1.35</u> (0.02)	<b>1177.69</b> (0.11)
Open Domain	<b>0.94</b> (0.09)	2675.39(6.14)	-	2732.16(0.73)	4.66(0.19)	<u>1498.12</u> (189.94)	<u>1.67</u> (0.02)	<b>1343.25</b> (0.18)
Overall	<b>0.93</b> (0.10)	2699.85(1.62)	-	2712.56(0.08)	4.53(0.16)	<u>1582.51</u> (207.25)	<u>1.53</u> (0.03)	<b>1271.69</b> (0.05)

sults in a substantial performance drop, indicating that extraction-based memories—once detached from raw dialogue context—rely heavily on semantic normalization processes such as coreference resolution and temporal alignment to maintain retrievability and consistency.

In contrast, enabling Knowledge Alignment for Episode Memory may even degrade performance. This observation suggests that when segmentation already achieves sufficient cognitive coherence, further semantic rewriting or fusion of raw dialogue can obscure implicit cues or fine-grained details. This result highlights a central design principle of HiMem: **semantic alignment should be memory-type aware rather than applied as a uniform preprocessing step.**

#### D.4 Memory Reconsolidation as a Key Mechanism for Long-Term Self-Evolution

Beyond static memory construction, experiments on Memory Self-Evolution further validate the importance of **Memory Reconsolidation** in sustained interactions. When Note Memory alone cannot support a query but Episode Memory provides sufficient evidence, the system performs targeted information extraction and conflict detection, writing missing knowledge back into Note Memory to correct and enrich its representation.

The effectiveness of this mechanism is evident at two levels. First, Note Memory exhibits a marked performance improvement after self-evolution is enabled, indicating substantive gains in knowledge coverage. Second, the corresponding improvement in overall performance demonstrates that **retrieval and memory updating should form a feedback loop rather than operate as independent processes.** This stands in contrast to memory systems

that only support additive or replacement-based updates, and more closely mirrors human cognition, where memory is reshaped through use.

#### D.5 Global Consistency Across Retrieval Modes and Efficiency Analyses

Finally, comparisons across retrieval modes and the Top- $K$  analysis provide system-level validation of the above design synergies. The complementary behaviors of hybrid and best-effort retrieval in terms of latency and token consumption reflect the flexibility afforded by hierarchical memory under an “abstraction-first, descend-when-necessary” strategy. Moreover, the observation that performance saturates at relatively small  $K$  values indicates that HiMem’s memory units possess high information density, avoiding inefficient compensation through expanded retrieval scopes.

#### D.6 Summary

In summary, HiMem’s advantages do not arise from larger models or longer context windows, but from a series of *design choices closely aligned with human long-term memory mechanisms*: cognitively consistent event segmentation, memory-type-aware semantic alignment, hierarchically complementary memory representations, and usage-driven Memory Reconsolidation. Together, these elements form a scalable, interpretable, and self-evolving long-term memory framework, providing robust support for sustained interaction and complex reasoning in long-horizon LLM agents.