

# NaturalMem: A Benchmark for Memory-Driven Dialogue in Large Language Models

Anonymous ACL submission

## Abstract

Most existing benchmarks for evaluating the memory of large language models (LLMs) rely on explicit recall-style question answering, where memory is directly queried and answered. However, such memory answering settings diverge from real world human-AI interaction, in which memory is rarely triggered explicitly and instead manifests implicitly by shaping dialogue generation. We introduce **NaturalMem**, a natural dialogue-based benchmark for evaluating **memory-driven dialogue**, where memory influences responses and speaking style without explicit recall prompts. **NaturalMem** constructs multi-turn dialogues based on fictional character prototypes, excluding identifiable names or show-specific entities. Each character is associated with personal facts, category information, and a target speaking style, and is evaluated across multiple dialogue sessions to assess fact retention and style consistency. Dialogue data is created through an LLM-assisted and human-curated pipeline. Experiments show that state-of-the-art agents still struggle to retain personal facts and maintain stable speaking styles in memory-driven dialogue settings. **NaturalMem** provides a realistic and diagnostic framework for evaluating memory in LLMs.

## 1 Introduction

Large language models (LLMs) are increasingly deployed as long-term conversational agents, personal assistants, and interactive systems that engage users over extended periods of time. In such settings, memory is a core capability: models are expected not only to process the immediate conversational context, but also to retain and consistently utilize user-specific information, personal facts, and behavioral patterns across interactions. Accordingly, evaluating the memory capabilities of LLMs and agent systems has emerged as an

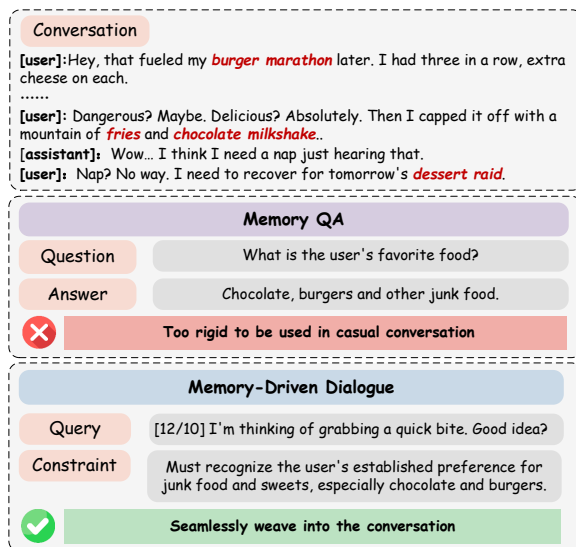


Figure 1: Transition from memory-based question answering to memory-driven dialogue.

important research direction, leading to a growing number of benchmarks that assess how models store, retrieve, and use past information (Gao et al., 2025; Hu et al., 2025; Wu et al., 2025; Zhang et al., 2025).

Most existing memory benchmarks, however, conceptualize memory primarily as an object of explicit querying. Typical evaluations adopt recall-style question answering, where models are directly asked to restate previously provided information (e.g., “What did the user say earlier?”) (Yang et al., 2018). While such **memory answering** paradigms are easy to quantify and provide clear evaluation targets, they rely on interaction patterns that are rare in natural human-AI conversations (Maharana et al., 2024; Wu et al., 2024; Zhong et al., 2024; Bai et al., 2024). In practice, users seldom probe a system’s memory explicitly. Instead, memory is expected to operate implicitly, shaping responses in a way that is consistent with prior interactions, as illustrated in Figure 1.

Table 1: Statistics of long-form conversational benchmarks from LongMemEval (Wu et al., 2024). These benchmarks differ in domain, scale, and interaction setting. DialSim is constructed from existing TV show transcripts, while NaturalMem consists of fictional yet naturalistic dialogues generated using real TV characters, specifically designed for controlled evaluation of long-term conversational memory.

Benchmark	Domain	Sess.	Ques.	Tokens / Conv.	Conversation Type
MSC(Xu et al., 2022)	Open-Domain	5k	–	1k	Human–Human
MemoryBank(Zhong et al., 2024)	Personal	300	194	5k	Human–AI
PerLTQA(Du et al., 2024)	Personal	4k	8,593	1M	Human–AI
LoCoMo(Maharana et al., 2024)	Personal	1k	7,512	10k	Human–Human
DialSim(Kim et al., 2024)	TV Shows	1k–2k	1M	350k	Human–Human
LongMemEval(Wu et al., 2024)	Personal	50k	500	115k, 1.5M	Human–AI
PersonaMem-V2(Jiang et al., 2025b)	Personal	–	5000	32k, 128k	Human–AI
<b>NaturalMem (ours)</b>	Fictional TV	27k	261	14k	Human–Human

This discrepancy reveals a fundamental gap between current benchmarks and realistic conversational use cases. In real-world dialogue, memory functions less as a mechanism for explicit recall and more as a behavioral driver of dialogue generation. A model demonstrates memory not by stating what it remembers, but by responding in ways that remain consistent with previously established facts, preferences, and interaction styles (Shuster et al., 2022). Importantly, this notion of memory-driven behavior differs from prompted or instruction-following personas. Rather than being repeatedly specified in the prompt, stylistic traits and personal characteristics are introduced implicitly through earlier interactions and must be retained and applied without explicit restatement.

To address this gap, we introduce **NaturalMem**, a benchmark designed to evaluate memory-driven dialogue in natural conversational settings. NaturalMem constructs dialogue scenarios based on characters drawn from television programs, whose coherent personal traits and stable speaking styles provide a realistic foundation for modeling long-term memory effects in dialogue. Each character is associated with a structured background specification, enabling systematic evaluation of both **factual consistency** and **stylistic stability** across interactions (Zhang et al., 2018). To prevent models from exploiting prior world knowledge, NaturalMem deliberately excludes character names and show-specific entities from the dialogue text, ensuring that successful performance requires genuine memory usage rather than surface-level recognition.

We evaluate a range of large language models and memory-augmented variants under this setting. Our results show that even models with strong explicit recall abilities often fail to consistently maintain personal facts and stable speaking styles when

memory must be applied implicitly in natural dialogue. These findings suggest that recall-centered benchmarks substantially underestimate the challenges of memory in realistic conversational interactions. NaturalMem is not intended to replace existing long-term memory benchmarks such as LoCoMo or LongMemEval, which focus on explicit recall or comprehension over extended dialogue histories. Instead, it targets a complementary and underexplored regime: whether memory can behaviorally and implicitly drive dialogue generation in the absence of explicit memory queries.

## 2 Related Work

### 2.1 Memory Modules for LLMs

Research on memory in language models spans from early neural architectures with explicit read–write memory to recent system-level memory augmentation for large language models. Early work introduced external memory mechanisms, such as Memory Networks (Weston et al., 2014), Neural Turing Machines (Graves et al., 2014), and Differentiable Neural Computers (Graves et al., 2016), with the goal of enabling models to store and retrieve information beyond the immediate input context.

Building on these foundations, recent systems incorporate practical memory components into LLM-based agents, including Mem0 (Chhikara et al., 2025), A-Mem (Xu et al., 2025), MemoryOS (Kang et al., 2025), and MemOS (Li et al., 2025). These approaches store user-specific facts or interaction histories and retrieve them to condition future generation, reflecting a broader trend toward treating memory as a first-class component in LLM applications rather than relying solely on extended context windows.

Despite their promise, such memory-augmented

138 systems are often evaluated using task-specific or  
139 ad hoc protocols. Consequently, it remains unclear  
140 whether the retrieved memory meaningfully influ-  
141 ences natural dialogue behavior, or merely serves  
142 as an external lookup mechanism, highlighting the  
143 need for systematic and realistic evaluation set-  
144 tings.

## 145 2.2 Memory Evaluation Benchmarks

146 A number of benchmarks have been proposed  
147 to evaluate memory and long-context capabilities  
148 of LLMs, most of which adopt recall-oriented  
149 paradigms. In these settings, memory is explic-  
150 itly queried through question answering, as exem-  
151 plified by MemoryBank (Zhong et al., 2024) and  
152 LongBench (Bai et al., 2024). More recent bench-  
153 marks such as LoCoMo (Maharana et al., 2024)  
154 and LongMemEval (Wu et al., 2024) extend this  
155 paradigm to longer and multi-session dialogues,  
156 yet still primarily rely on explicit probing.

157 Beyond evaluation protocols, existing long-  
158 form conversational memory benchmarks also vary  
159 widely in their domains, scales, and interaction set-  
160 tings. As summarized in Table 1, prior datasets  
161 differ in whether conversations are Human–Human  
162 or Human–AI and the typical conversation length.  
163 Consequently, these benchmarks predominantly  
164 measure memory answering, rather than whether  
165 memory implicitly shapes dialogue generation in  
166 realistic conversational settings.

## 167 2.3 Implicit Memory and Style Consistency

168 Beyond factual recall, prior work has empha-  
169 sized persona and stylistic consistency as impor-  
170 tant aspects of dialogue systems. Persona-based  
171 benchmarks such as Persona-Chat (Zhang et al.,  
172 2018) show that maintaining stable personal at-  
173 tributes improves conversational coherence. More  
174 recent efforts further explore related dimensions:  
175 LikeBench (Awsafur Rahman et al., 2025) evalu-  
176 ates likability and adaptability in multi-session per-  
177 sonalized interactions, while PersonaMem (Jiang  
178 et al., 2025b,a) focuses on modeling implicit per-  
179 sonas and evolving user profiles over time.

180 These studies underscore that memory in dia-  
181 logue systems is often expressed implicitly, through  
182 consistent behavior and linguistic style rather than  
183 explicit recollection. However, existing bench-  
184 marks typically evaluate persona or style in iso-  
185 lation, without jointly assessing factual consistency  
186 and stylistic continuity in natural, multi-session  
187 dialogue.

NaturalMem differs from prompted or  
instruction-following persona settings commonly  
studied in prior work. Prompted personas rely on  
explicit conditioning, where stylistic attributes are  
repeatedly specified or reinforced in the prompt.  
In contrast, the style consistency evaluated in  
NaturalMem is memory-driven: stylistic traits are  
introduced implicitly through earlier interactions  
and must be retained and applied in later sessions  
without restatement.

Maintaining such behavior therefore requires  
models to utilize long-term contextual information  
across sessions rather than responding to local in-  
structions. In this sense, style consistency in Nat-  
uralMem constitutes a form of implicit memory,  
manifested behaviorally through stable linguistic  
patterns, tone, and interactional preferences over  
time, even in the presence of distractors.

## 3 Reverse-Guided Dataset Construction

Fully automatic dataset construction using large  
language models alone (Jiang et al., 2025b), corre-  
sponding to LLM draft generation without human  
intervention (Figure 2, Stage 2), is often insuffi-  
cient for natural conversational memory evaluation.  
In practice, LLM-generated drafts frequently vio-  
late de-identification constraints, exhibit repetitive  
question patterns, and fail to support deep, dis-  
tributed memory dependencies required by down-  
stream evaluation (Figure 2, Stage 4).

To address these limitations, we propose  
**Reverse-Guided Dataset Construction (RGDC)**,  
a human-intensive (Long et al., 2024; Maharana  
et al., 2024) and iterative pipeline that treats LLM  
outputs as intermediate drafts rather than final arti-  
facts. As illustrated in Figure 2, RGDC integrates  
structured human curation and feedback (Figure 2,  
Stage 3) to iteratively correct generation failures  
and guide constrained regeneration.

### 3.1 LLM-Based Draft Dialogue Generation

As shown in Figure 2 (Stage 2), dialogues are  
first generated by an LLM conditioned on abstract  
role descriptions and multi-session conversational  
contexts. Explicit de-identification constraints are  
imposed to discourage character names and other  
identity-revealing references, ensuring that correct  
answers must be supported by prior conversational  
evidence rather than external knowledge.

Despite these constraints, drafts produced at this  
stage often contain identity leakage, stereotypical

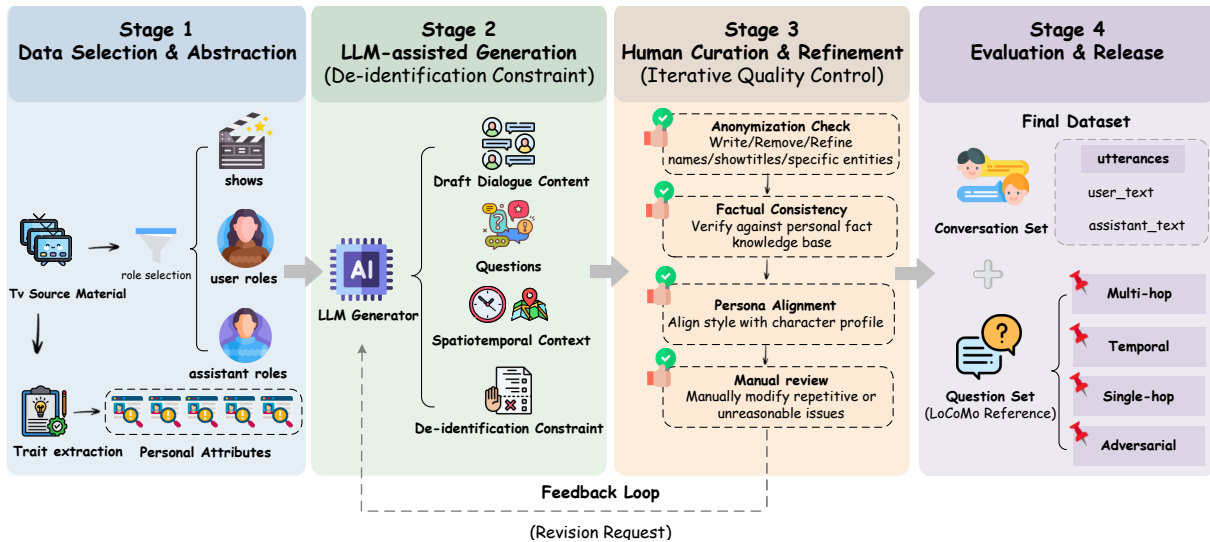


Figure 2: Reverse-Guided Dataset Construction (RGDC), an iterative human-in-the-loop pipeline. The process consists of four stages: (S1) data selection and abstraction, (S2) LLM-assisted draft generation under de-identification constraints, (S3) human curation and refinement with structured quality control, and (S4) evaluation-oriented dataset assembly. Common failure modes in Stage 2 are addressed through iterative feedback and constrained revision in Stage 3, forming a closed loop that targets the reasoning requirements of Stage 4.

expressions, and shallow memory cues. All generated dialogues are therefore treated as provisional drafts and forwarded to human curation (Figure 2, Stage 3).

### 3.2 Human Curation and Structural Enforcement

To correct systematic failures from LLM drafting (Figure 2, Stage 2), all dialogues undergo manual inspection during human curation (Figure 2, Stage 3). This step focuses on enforcing structural constraints that are difficult to guarantee through prompting alone, particularly the removal of identity leakage and shortcut signals that could enable reliance on pretraining knowledge.

Edits are performed through targeted manual rewriting rather than automatic deletion, as naive removal often disrupts discourse coherence (Wu et al., 2024).

### 3.3 Reverse-Guided Question Construction

Evaluation questions are explicitly decoupled from initial dialogue generation (Figure 2) and constructed in a reverse-guided manner. Questions generated directly by LLMs tend to be repetitive and biased toward simple memory patterns.

Human annotators therefore design evaluation questions first (Figure 2, Stage 4), targeting diverse reasoning requirements, including multi-hop, temporal, and adversarial memory dependencies.

Dialogue content is subsequently revised or regenerated so that answers are naturally grounded in prior conversations, forming an explicit feedback loop (Figure 2).

### 3.4 Iterative Human-in-the-Loop Refinement

Overall, RGDC operates as an iterative human-in-the-loop process rather than a linear pipeline, as explicitly illustrated by the feedback loop in Figure 2. Human verification and editing are carried out by multiple annotators drawn from the same research group, all of whom are familiar with the construction guidelines and independently review the generated dialogues. Revisions are cross-checked by at least one additional annotator to reduce individual bias and ensure consistency.

Dialogues often undergo multiple rounds of inspection, manual editing, constrained regeneration, and re-verification across Stage 2 and Stage 3. Correcting one issue (e.g., removing an identity cue) may introduce others (e.g., reduced coherence or loss of implicitness), requiring further iteration. Question sets are also repeatedly reviewed to eliminate redundancy; when duplicate or overly similar questions are identified, annotators manually revise or replace them, followed by corresponding dialogue adjustments. This iterative refinement continues until dialogues satisfy constraints on naturalness, memory dependency, and robustness, or are discarded.

## 4 Evaluation Benchmark

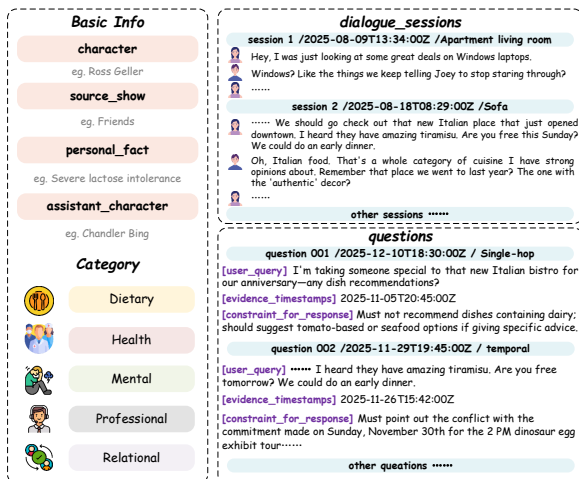


Figure 3: Schema of our constructed dataset. Each instance comprises a user character persona derived from a TV show, an assistant character style, multiple anonymized dialogue sessions, and a set of evaluation questions with grounded factual constraints.

**Dataset overview.** Based on the multi-session dialogues constructed in §3, we introduce an evaluation benchmark designed to measure long-term conversational memory under identity-anonymized settings. The dataset comprises **73** unique source\_show, **90** distinct character personas, and **89** assistant\_character styles. In total, it contains **2,774** dialogue\_sessions with **43,492** utterances, amounting to approximately **1.3M** tokens. The structure of a single instance is illustrated in Figure 3. We further annotate **261** evaluation questions spanning diverse long-term memory reasoning requirements.

Crucially, all dialogues are identity-anonymized: character names and show-specific entities are removed from the text. This design ensures that successful performance cannot be achieved through entity recognition or prior world knowledge, but instead requires models to retain and utilize dialogue-grounded memory across sessions.

**Personal fact categories.** To characterize the content of conversational memory, we categorize personal facts appearing in dialogues into five types: Dietary, Health, Mental, Professional, and Relational. These categories reflect common forms of personal information naturally exchanged in everyday conversations and are designed to capture diverse memory demands.

Different categories place distinct requirements

on long-term memory. For example, professional routines often involve stable, recurring facts, while mental or relational states may evolve gradually across sessions. By structuring memory targets along these dimensions, NaturalMem enables more fine-grained analysis of how models retain and utilize different types of personal information over time.

**Question categories.** Following prior work on multi-session memory evaluation (Maharana et al., 2024), we categorize evaluation questions into four types: (1) **Single-hop** questions depend on a single memory point, typically localized within one session; (2) **Multi-hop** questions require integrating multiple historical facts distributed across sessions; (3) **Temporal** questions involve time ordering, state transitions, or before/after relations that cannot be resolved from a single isolated mention; (4) **Adversarial** questions include misleading, conflicting, or underspecified cues, testing whether models can robustly preserve grounded facts and reject unsupported inferences.

This taxonomy is designed to reflect realistic conversational memory usage while preventing shortcut answering via character identity, show-specific knowledge, or external world facts. Unlike some prior benchmarks, we do not introduce a separate open-domain question category. This choice is intentional. In our identity-anonymized setting, every evaluation question already requires open-ended natural language generation grounded in dialogue history, rather than selecting or extracting a short factual span. Models must integrate recalled information with open-domain reasoning to produce a coherent response that satisfies the specified constraints. Consequently, open-domain reasoning is not treated as an isolated category, but is instead an inherent property of all question types in NaturalMem.

### 4.1 Task 1: Factual Memory Consistency

A conversational agent is expected to recall dialogue-grounded personal facts introduced implicitly over long histories and to maintain factual consistency in subsequent sessions. We formulate **Fact** as an accuracy-based evaluation. For each query, the agent generates a response that must satisfy a set of factual requirements encoded in constraint\_for\_response. A response is marked as correct if and only if it satisfies all specified constraints and does not contradict any

Table 2: Performance of base LLMs across different question categories and overall. Fact is reported as accuracy (%), Style as mean Likert score (1–5), and Latency as average wall-clock generation time (seconds) per query. All evaluations are conducted using DeepSeek-V3.2 as the unified LLM-as-a-Judge.

Model	Single-hop		Multi-hop		Temporal		Adversarial		Overall		Latency (s)
	Fact	Style	Fact	Style	Fact	Style	Fact	Style	Fact	Style	
DeepSeek-V3.2	75.6	4.07	64.9	4.23	<b>39.0</b>	4.06	54.5	4.01	55.6	4.08	3.37
Doubao-1.5-Pro-32k	77.8	4.13	<b>75.4</b>	4.21	30.5	4.11	<b>68.8</b>	4.09	<b>59.8</b>	4.13	6.07
Qwen3-235b-a22b	77.3	4.14	65.5	<b>4.55</b>	29.1	4.35	58.7	<b>4.40</b>	54.2	4.37	11.67
Qwen3-32b	79.5	4.21	56.4	4.53	30.4	4.28	50.7	4.15	50.6	4.28	6.99
Qwen3-30b-a3b	79.5	<b>4.50</b>	70.9	<b>4.55</b>	27.8	<b>4.38</b>	52.0	4.37	53.4	<b>4.44</b>	4.55
GPT-4o-mini	68.9	3.49	59.6	3.91	24.4	3.98	41.6	3.88	44.8	3.85	3.95
GPT-4.1-mini	<b>82.2</b>	3.49	66.7	3.79	<b>39.0</b>	3.68	42.9	3.92	53.6	3.74	<b>2.95</b>

dialogue-grounded facts; otherwise, it is marked incorrect. We report overall **Accuracy**, as well as accuracy broken down by question category.

**LLM-as-a-Judge.** We adopt an LLM-as-a-Judge protocol to determine factual correctness. Automatic n-gram overlap metrics such as BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and token-level F1 are ill-suited for this task, as valid answers may differ substantially in surface form while preserving identical factual content, and high lexical overlap may still conceal subtle contradictions. The judge model is instructed to verify constraint satisfaction based on the dialogue history and to output a binary correctness decision, along with a brief rationale for auditing purposes.

## 4.2 Task 2: Assistant Style Consistency

Beyond factual memory, natural long-term interaction also requires stable stylistic behavior. In the **Style** task, the assistant is expected to respond in a manner consistent with a target assistant style described by `assistant_character`, across multiple sessions and in the presence of distractors. Notably, the style specification is introduced implicitly through earlier interactions and is not restated in the evaluation prompt.

We evaluate style consistency using an LLM-as-a-Judge protocol that compares the generated response against the target style specification.

Scoring (1–5). Each response is rated on a five-point Likert scale (Likert, 1932): **1** = clearly off-style with strong drift; **2** = mostly off-style with occasional matching cues; **3** = partially consistent but unstable or generic; **4** = largely consistent with minor deviations; **5** = highly consistent and characterful throughout.

**LLM-as-a-Judge.** We use LLM-based evaluation because stylistic fidelity is inherently non-

lexical and cannot be reliably captured by overlap-based metrics. BLEU, ROUGE, and F1 reward surface copying and penalize legitimate stylistic variation, while failing to measure higher-level signals such as tone, humor, formality, or verbosity specified in `assistant_character`. The judge is instructed to assess stylistic faithfulness independently of factual correctness.

## 4.3 Task 3: Answering Latency under Memory Retrieval

Memory-augmented agents may incur additional overhead due to retrieval, re-ranking, and prompt construction. We therefore measure end-to-end response time per query under a standardized inference setting. Latency statistics, including mean values, are reported to characterize efficiency–accuracy trade-offs across models.

## 5 Experiment Setup

All experiments are conducted under a unified text-only setting. Models are evaluated on anonymized multi-session dialogues without access to character identities, show-specific entities, or any external metadata. All systems are tested on the same dialogue histories and evaluation questions described in §4. To ensure comparability across models and tasks, all evaluations are standardized using a single fixed evaluation model.

**Models.** We evaluate both base large language models and memory-augmented agents to assess long-term conversational memory under diverse modeling paradigms. As base models, we include a representative set of state-of-the-art LLMs spanning different architectures, training regimes, and parameter scales: **DeepSeek-V3.2**, **Doubao-1.5-Pro-32k-250115**, **GPT-4o-mini**, **GPT-4.1-mini**, as well as three models from the **Qwen3** family:

Table 3: Performance of memory-augmented agents across different question categories and overall. Fact is reported as accuracy (%), Style as mean Likert score (1–5), and Latency as average wall-clock generation time (seconds) per query. All evaluations are conducted using DeepSeek-V3.2 as the unified LLM-as-a-Judge.

Method	Single-hop		Multi-hop		Temporal		Adversarial		Overall		Latency (s)
	Fact	Style	Fact	Style	Fact	Style	Fact	Style	Fact	Style	
<b>GPT-4o-mini</b>											
Mem0	40.0	2.38	47.4	2.33	18.3	2.46	6.5	2.48	24.9	2.43	4.76
A-Mem	48.9	2.93	63.2	2.96	<b>32.9</b>	3.24	16.9	3.08	37.5	3.08	3.61
MemoryOS	64.4	3.33	<b>70.2</b>	3.39	30.5	3.24	<b>51.9</b>	3.38	<b>51.3</b>	3.33	<b>2.39</b>
MemOS	46.7	3.07	59.6	3.09	11.0	3.15	3.9	3.34	25.7	3.17	3.58
Full-context	<b>68.9</b>	<b>3.49</b>	59.6	<b>3.91</b>	24.4	<b>3.98</b>	41.6	<b>3.88</b>	44.8	<b>3.85</b>	3.95
<b>GPT-4.1-mini</b>											
Mem0	37.8	2.36	42.1	2.21	23.2	2.52	6.5	2.62	24.9	2.46	4.38
A-Mem	64.4	3.00	71.9	3.16	36.6	3.44	28.6	3.40	46.7	3.29	2.91
MemoryOS	64.4	<b>3.82</b>	<b>73.7</b>	3.77	31.7	3.51	<b>50.6</b>	3.79	52.1	3.70	<b>2.23</b>
MemOS	48.9	3.29	70.2	3.44	13.4	3.51	6.5	3.65	29.9	3.50	4.36
Full-context	<b>82.2</b>	3.49	66.7	<b>3.79</b>	<b>39.0</b>	<b>3.68</b>	42.9	<b>3.92</b>	<b>53.6</b>	<b>3.74</b>	2.95

**Qwen3-235B-A22B, Qwen3-32B, and Qwen3-30B-A3B.** This selection covers both dense and mixture-of-experts architectures, and spans a wide range of model capacities, enabling a more comprehensive evaluation of memory behavior across scales.

To assess the effectiveness of explicit memory mechanisms beyond extended context, we additionally evaluate representative memory-augmented agents, including **Mem0**, **A-Mem**, **MemoryOS**, and **MemOS**. All memory-augmented agents are configured to store and retrieve dialogue histories according to their recommended settings. For fair comparison, identical dialogue inputs, prompts, and evaluation protocols are used across base LLMs and memory-augmented agents.

**Evaluation Protocol.** Factual memory consistency and assistant style consistency are evaluated using the LLM-as-a-Judge protocol defined in §4. Judgments for factual correctness and stylistic fidelity are performed independently. All models generate responses under the same prompting and decoding settings to control for generation variability.

**Latency Measurement.** To characterize efficiency trade-offs, we measure end-to-end response latency for all models when answering benchmark queries. Latency is recorded as wall-clock time per query under a standardized inference environment. For memory-augmented agents, the reported latency includes memory retrieval, re-ranking, and prompt construction overhead. We report average latency as well as percentile statistics to re-

flect practical deployment considerations. Detailed hardware configurations and runtime settings are provided in the Appendix A.1.

**Evaluation Model.** We use **DeepSeek-V3.2** as the primary evaluation model for all LLM-as-a-Judge judgments, including factual correctness verification and style scoring. Using a single evaluator ensures consistency and comparability across models and tasks.

To further verify robustness and diversity, we additionally employ **GPT-4o-mini** as an alternative judge model and perform manual spot-checks on randomly sampled responses. The overall trends and conclusions remain consistent across both automatic judges and human verification; detailed comparisons are provided in Appendix B.

## 6 Results and Analysis

We present the experimental results on NaturalMem from four perspectives: model scale effects, factual memory versus stylistic consistency, the behavior of memory-augmented agents, and efficiency–accuracy trade-offs. Additional analysis and discussion are provided in Appendix A.2.

### 6.1 Effect of Model Scale

Table 2 reports the performance of base LLMs on NaturalMem. Overall, we observe that larger or more capable models do not consistently achieve higher accuracy on this benchmark. In particular, improvements in single-hop factual questions do not reliably transfer to multi-hop, temporal, or adversarial settings.

This result contrasts with trends commonly observed in explicit recall or short-context reasoning benchmarks, indicating that NaturalMem captures a distinct aspect of conversational memory beyond raw model capacity.

## 6.2 Factual Memory vs. Style Consistency

As shown in Table 2, factual accuracy varies substantially across question categories, with the largest drops occurring in temporal and adversarial questions. In contrast, style consistency scores are relatively stable across models and categories.

This discrepancy suggests that factual consistency is more sensitive to long-term memory degradation than stylistic behavior.

Table 4: Accuracy (%) comparison between LoCoMo and NaturalMem.  $\Delta$  denotes the accuracy difference.

Model	LoCoMo	NaturalMem	$\Delta$
Mem0	66.3	24.9	41.4
A-Mem	62.1	37.5	24.6
MemoryOS	61.6	51.3	10.3
MemOS	73.3	25.7	47.6
Full-context	80.6	44.8	35.8

## 6.3 Memory-Augmented Agents

Results for memory-augmented agents are summarized in Table 3. Compared to base models, several agents achieve higher factual accuracy, demonstrating the benefit of explicit memory retrieval, but often at the cost of lower style consistency relative to the full-context baseline.

These results further suggest that memory selection itself may be a source of error, as retrieved content does not always align with the requirements of the current query. Among the evaluated systems, MemoryOS exhibits a more balanced performance across factual accuracy and style consistency, highlighting trade-offs in current agent-based memory designs.

## 6.4 Full-Context Baselines

Despite the absence of explicit memory retrieval, the full-context baseline achieves competitive and often higher factual accuracy than several memory-augmented agents (Table 3), indicating that imperfect retrieval and memory selection can partially offset the benefits of external memory mechanisms.

Notably, full-context models maintain substantially higher style consistency than most agent-based systems, suggesting that current agents pri-

marily focus on fact-oriented memory extraction while stylistic cues are less effectively preserved.

However, full-context inference incurs higher computational costs, motivating the use of memory-augmented approaches in practical deployments.

## 6.5 Latency Analysis

Latency statistics are reported in Tables 2 and 3. For base models and full-context inference, latency increases with input length and is dominated by prefill and generation costs. Memory-augmented agents introduce additional overhead from memory retrieval and prompt construction, resulting in different efficiency–accuracy trade-offs.

## 6.6 Comparison with LoCoMo

Table 4 compares performance on LoCoMo and NaturalMem. All evaluated systems perform substantially worse on NaturalMem, confirming that it poses a more challenging evaluation setting.

Unlike LoCoMo, which explicitly probes memory through recall-style questions, NaturalMem evaluates whether memory can implicitly guide response generation in natural dialogue, as illustrated in Figure 1.

## 7 Conclusion and Future Work

We introduce **NaturalMem**, a benchmark for evaluating *memory-driven dialogue*, where memory implicitly guides factual consistency and speaking style rather than being explicitly queried. Our results show that both base LLMs and memory-augmented agents struggle under this setting, despite strong performance on recall-centric benchmarks, revealing limitations of existing memory evaluation protocols.

### Several directions remain for future work.

First, our current agent evaluation assumes that memory is always retrieved. In realistic deployments, agents must additionally decide when memory should be stored and whether retrieval is necessary for a given query. The ability to selectively trigger memory retrieval may not only improve memory utilization, but also reduce response latency by avoiding unnecessary retrieval and prompt construction. Second, extending NaturalMem to multimodal interactions may enable richer assessments of personalized and long-term consistency, moving toward a more comprehensive form of personalized Turing-style evaluation.

## 597 Limitations

598 NaturalMem has several **limitations**. First, the  
599 number of evaluation questions is relatively small,  
600 as question construction requires substantial hu-  
601 man involvement to ensure naturalness and robust  
602 memory dependencies. Second, our agent evalua-  
603 tion does not model autonomous decisions about  
604 when to store or retrieve memory, focusing instead  
605 on memory utilization given retrieved content. Fi-  
606 nally, while we employ LLM-as-a-Judge evaluation  
607 with manual spot-checking, automatic evaluation  
608 of memory-driven dialogue remains an open chal-  
609 lenge.

## 610 References

611 Md Awsafur Rahman, Adam Gabrys, Doug Kang,  
612 Jingjing Sun, Tian Tan, and Ashwin Chandramouli.  
613 2025. Likebench: Evaluating subjective likability  
614 in llms for personalization. *arXiv e-prints*, pages  
615 arXiv-2512.

616 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu,  
617 Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao  
618 Liu, Aohan Zeng, Lei Hou, and 1 others. 2024. Long-  
619 bench: A bilingual, multitask benchmark for long  
620 context understanding. In *Proceedings of the 62nd*  
621 *annual meeting of the association for computational*  
622 *linguistics (volume 1: Long papers)*, pages 3119–  
623 3137.

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

628 Yiming Du, Hongru Wang, Zhengyi Zhao, Bin Liang,  
629 Baojun Wang, Wanjun Zhong, Zezhong Wang, and  
630 Kam-Fai Wong. 2024. Perltqa: A personal long-term  
631 memory dataset for memory classification, retrieval,  
632 and synthesis in question answering. *arXiv preprint*  
633 *arXiv:2402.16288*.

634 Huan-ang Gao, Jiayi Geng, Wenyue Hua, Mengkang Hu,  
635 Xinzhe Juan, Hongzhang Liu, Shilong Liu, Jiahao  
636 Qiu, Xuan Qi, Yiran Wu, and 1 others. 2025. A  
637 survey of self-evolving agents: On path to artificial  
638 super intelligence. *arXiv preprint arXiv:2507.21046*.

639 Alex Graves, Greg Wayne, and Ivo Danihelka.  
640 2014. Neural turing machines. *arXiv preprint*  
641 *arXiv:1410.5401*.

642 Alex Graves, Greg Wayne, Malcolm Reynolds,  
643 Tim Harley, Ivo Danihelka, Agnieszka Grabska-  
644 Barwińska, Sergio Gómez Colmenarejo, Edward  
645 Grefenstette, Tiago Ramalho, John Agapiou, and  
646 1 others. 2016. Hybrid computing using a neural  
647 network with dynamic external memory. *Nature*,  
648 538(7626):471–476.

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 J  
Taylor, and Dan Roth. 2025a. Know me, respond to  
me: Benchmarking llms for dynamic user profiling  
and personalized responses at scale. *arXiv preprint*  
*arXiv:2504.14225*.

Bowen Jiang, Yuan Yuan, Maohao Shen, Zhuoqun Hao,  
Zhangchen Xu, Zichen Chen, Ziyi Liu, Anvesh Rao  
Vijjini, Jiashu He, Hanchao Yu, Radha Poovendran,  
Gregory Wornell, Lyle Ungar, Dan Roth, Sihao Chen,  
and Camillo Jose Taylor. 2025b. [Personamem-v2: Towards personalized intelligence via learning implicit user personas and agentic memory](#). *Preprint*,  
arXiv:2512.06688.

Jiazheng Kang, Mingming Ji, Zhe Zhao, and Ting  
Bai. 2025. Memory os of ai agent. *arXiv preprint*  
*arXiv:2506.06326*.

Jiho Kim, Woosog Chay, Hyeonji Hwang, Daeun  
Kyung, Hyunseung Chung, Eunbyeol Cho, Yohan  
Jo, and Edward Choi. 2024. Dialsim: A real-time  
simulator for evaluating long-term dialogue under-  
standing of conversational agents. *arXiv e-prints*,  
pages arXiv-2406.

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

Rensis Likert. 1932. A technique for the measurement  
of attitudes. *Archives of psychology*.

Chin-Yew Lin. 2004. Rouge: A package for automatic  
evaluation of summaries. In *Text summarization*  
*branches out*, pages 74–81.

Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao  
Ding, Gang Chen, and Haobo Wang. 2024. On llms-  
driven synthetic data generation, curation, and evalu-  
ation: A survey. *arXiv preprint arXiv:2406.15126*.

Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov,  
Mohit Bansal, Francesco Barbieri, and Yuwei  
Fang. 2024. Evaluating very long-term conver-  
sational memory of llm agents. *arXiv preprint*  
*arXiv:2402.17753*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-  
Jing Zhu. 2002. Bleu: a method for automatic evalu-  
ation of machine translation. In *Proceedings of the*  
*40th annual meeting of the Association for Computa-*  
*tional Linguistics*, pages 311–318.

Kurt Shuster, Jing Xu, Mojtaba Komeili, Da Ju,  
Eric Michael Smith, Stephen Roller, Megan Ung,

649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702

703	Moya Chen, Kushal Arora, Joshua Lane, and 1 others. 2022. Blenderbot 3: a deployed conversational agent that continually learns to responsibly engage. <i>arXiv preprint arXiv:2208.03188</i> .	<b>Section B</b>	reports additional validation of the LLM-as-a-Judge protocol, including alternative judge models and human verification.	757
704				758
705				759
706				
707	Jason Weston, Sumit Chopra, and Antoine Bordes. 2014. Memory networks. <i>arXiv preprint arXiv:1410.3916</i> .	<b>Section C</b>	provides the prompt templates used in experiments.	760
708				761
709	Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang, Kai-Wei Chang, and Dong Yu. 2024. Longmemeval: Benchmarking chat assistants on long-term interactive memory. <i>Preprint</i> , arXiv:2410.10813.	<b>A</b>	<b>Experiment Details and Analysis</b>	762
710		<b>A.1</b>	<b>Evaluation Procedures</b>	763
711			We detail the evaluation procedures used for base models and memory-augmented agents to ensure reproducibility and fair comparison.	764
712			For base LLMs and the full-context baseline, the entire dialogue history is provided as input, with the final turn being the user query. Models are instructed to directly generate a response without any additional task-specific prompts or instructions.	765
713	Yaxiong Wu, Sheng Liang, Chen Zhang, Yichao Wang, Yongyue Zhang, Huifeng Guo, Ruiming Tang, and Yong Liu. 2025. From human memory to ai memory: A survey on memory mechanisms in the era of llms. <i>arXiv preprint arXiv:2504.15965</i> .		For memory-augmented agents, dialogue histories are first processed to construct their internal memory representations according to each agent’s recommended configuration. At inference time, retrieved memory content is concatenated with the user query using a minimal prompt template, and the model generates a response conditioned on the retrieved memory.	766
714			All models and agents are evaluated via API calls under comparable network conditions. Since all systems are invoked through API-based interfaces, network latency is consistent across models and can be safely ignored. Reported latency therefore primarily reflects model inference and memory retrieval overhead, ensuring the reliability of latency comparisons.	767
715		<b>A.2</b>	<b>Additional Analysis of Results (Tables 2–3)</b>	768
716			This subsection provides a detailed, table-grounded analysis to complement the high-level observations in §6.	769
717			<b>(1) Category-wise factual difficulty: Temporal and adversarial are the bottlenecks.</b> Table 2 shows a consistent pattern across base LLMs: <i>Single-hop</i> factual accuracy is relatively high (e.g., 68.9–82.2), <i>Multi-hop</i> drops moderately (56.4–75.4), while <i>Temporal</i> exhibits the largest degradation (24.4–39.0), and <i>Adversarial</i> remains challenging (41.6–68.8). This explains why gains on single-hop do not reliably transfer to more realistic settings: NaturalMem’s hardest categories require maintaining state transitions and rejecting misleading cues, which are precisely where models fail	770
718	Jing Xu, Arthur Szlam, and Jason Weston. 2022. Beyond goldfish memory: Long-term open-domain conversation. In <i>Proceedings of the 60th annual meeting of the association for computational linguistics (volume 1: long papers)</i> , pages 5180–5197.			771
719				772
720				773
721				774
722				775
723	Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. 2025. A-mem: Agentic memory for llm agents. <i>arXiv preprint arXiv:2502.12110</i> .			776
724				777
725				778
726				779
727	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.			780
728				781
729				782
730				783
731				784
732				785
733				786
734				787
735	Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? <i>arXiv preprint arXiv:1801.07243</i> .			788
736				789
737				790
738				791
739	Zeyu Zhang, Quanyu Dai, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. 2025. A survey on the memory mechanism of large language model-based agents. <i>ACM Transactions on Information Systems</i> , 43(6):1–47.			792
740				793
741				794
742				795
743				796
744	Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. 2024. Memorybank: Enhancing large language models with long-term memory. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19724–19731.			797
745				798
746				799
747				800
748				801
749	<b>Appendix Overview</b>			802
750	The appendix provides supplementary details to support reproducibility and further analysis of NaturalMem.			803
751				804
752				
753	<b>Section A</b> describes the experimental setup, evaluation procedures, and additional analyses of memory behavior for both base models and memory-augmented agents.			
754				
755				
756				

805 most.

806 **(2) Model scale is not predictive under memory-**  
807 **driven dialogue.** Despite large model capacity,  
808 Qwen3-235B-A22B achieves 54.2 overall factual  
809 accuracy, comparable to DeepSeek-V3.2 (55.6) and  
810 below Doubao-1.5-Pro-32k (59.8), while incurring  
811 much higher latency (11.67s vs. 3.37s and 6.07s,  
812 respectively). This non-monotonicity suggests that  
813 memory-driven dialogue performance depends less  
814 on raw capacity and more on whether the model  
815 can robustly *select* and *apply* distributed dialogue  
816 evidence, especially for temporal/adversarial con-  
817 straints.

818 **(3) Fact vs. style: style scores are comparatively**  
819 **stable, but do not imply strong factual memory.**  
820 Across base LLMs, style consistency is relatively  
821 high and clustered (roughly 3.74–4.55), even when  
822 factual accuracy differs substantially (44.8–59.8  
823 overall). For example, GPT-4o-mini attains 3.85  
824 style with 44.8 fact, whereas Doubao reaches 4.13  
825 style with 59.8 fact. This indicates that producing  
826 a style-consistent response is easier than satisfying  
827 strict factual constraints: a response can be stylistically  
828 plausible while still violating a constraint  
829 (e.g., temporal ordering, exclusions).

830 **(4) Agent systems: factual gains come with a**  
831 **large style drop, especially versus full-context.**  
832 Table 3 reveals a clear trade-off for memory-  
833 augmented agents. Under GPT-4o-mini, MemoryOS  
834 improves overall factual accuracy to 51.3, sur-  
835 passing the full-context baseline (44.8), while also  
836 reducing latency (2.39s vs. 3.95s). However, its  
837 style score drops to 3.33 compared to full-context  
838 3.85. A similar pattern holds for A-Mem (fact  
839 37.5; style 3.08) and Mem0 (fact 24.9; style 2.43),  
840 both far below full-context style. This supports  
841 the hypothesis that current agent memory pipelines  
842 are primarily *fact-centric*: they retrieve and inject  
843 factual summaries well enough to help constraint  
844 satisfaction, but do not preserve the implicit stylistic  
845 cues accumulated across sessions, leading to  
846 style drift.

847 **(5) Memory selection is a measurable error**  
848 **source.** The agent results also suggest that re-  
849 trieval *relevance/selection* is a key limitation. For  
850 instance, MemOS under GPT-4o-mini attains rea-  
851 sonable style (3.17) but very low factual accuracy  
852 (25.7), with catastrophic performance on adver-  
853 sarial (3.9) and temporal (11.0). Such patterns are  
854 consistent with retrieving incomplete or misaligned

855 memories: the generation remains fluent (and some-  
856 times stylistically acceptable) while failing hard on  
857 constraint satisfaction.

858 **(6) Full-context remains a strong upper bound**  
859 **for style, and often for facts.** Full-context gen-  
860 erally provides higher style consistency than agent  
861 systems. For GPT-4o-mini, full-context style is  
862 3.85, above MemoryOS 3.33 and A-Mem 3.08; for  
863 GPT-4.1-mini, full-context is 3.74, close to Memo-  
864 ryOS 3.70 and above A-Mem 3.29. On factual ac-  
865 curacy, full-context is also competitive: it reaches  
866 53.6 on GPT-4.1-mini, slightly higher than Memo-  
867 ryOS 52.1, indicating that retrieval benefits can be  
868 offset by imperfect memory selection and compres-  
869 sion.

870 **(7) Latency trade-offs: retrieval can be faster**  
871 **than full-context, but not always.** Latency pro-  
872 files differ substantially. On GPT-4o-mini, Memo-  
873 ryOS reduces latency (2.39s) below full-context  
874 (3.95s) while improving facts, but Mem0 increases  
875 latency (4.76s) while performing much worse  
876 (24.9). On GPT-4.1-mini, MemoryOS is again  
877 faster (2.23s) than full-context (2.95s). These re-  
878 sults suggest that retrieval-based pipelines can re-  
879 duce end-to-end latency when memory selection  
880 is accurate and retrieval overhead is controlled;  
881 otherwise, retrieval adds cost without improving  
882 correctness. This motivates future work on *adap-*  
883 *tive retrieval triggering* (deciding when retrieval  
884 is necessary), which could further reduce latency  
885 by avoiding unnecessary retrieval and prompt con-  
886 struction.

## 887 B Judge Model Robustness

### 888 B.1 Robustness to the Choice of 889 LLM-as-a-Judge

890 In addition to using **DeepSeek-V3.2** as the pri-  
891 mary LLM-as-a-Judge throughout our main exper-  
892 iments, we evaluate the robustness of factual cor-  
893 rectness assessment by replacing the judge model  
894 with **GPT-4o-mini**. Crucially, all agent responses  
895 are kept identical across settings; only the evaluator  
896 is changed.

897 Table 5 reports the overall factual consistency  
898 of memory-augmented agents under the two eval-  
899 uators. Although absolute accuracy values differ  
900 across judges, the relative performance trends re-  
901 main stable. In particular, **MemoryOS** and **Full-**  
902 **context** consistently achieve the highest factual  
903 accuracy under both judges, whereas lightweight

Table 5: Overall factual consistency (%) of memory-augmented agents under different LLM-as-a-Judges. Agent responses are identical across settings; only the evaluator is changed. This table is included in the appendix as a robustness check.

Method	DeepSeek-V3.2 Judge	GPT-4o-mini Judge
Mem0	24.9	29.5
A-Mem	46.7	51.7
MemoryOS	52.1	<b>59.4</b>
MemOS	29.9	30.7
Full-context	<b>53.6</b>	56.7

memory mechanisms such as **Mem0** and **MemOS** remain substantially weaker. This indicates that our main conclusions are not sensitive to the choice of LLM-based evaluator.

We further observe that GPT-4o-mini generally assigns higher factual correctness scores than DeepSeek-V3.2. Since the evaluated agent responses are identical, this discrepancy reflects differences in judgment strictness and calibration rather than changes in agent behavior. Such systematic shifts are expected when switching between large language models with different evaluation granularity and error tolerance, and do not affect comparative conclusions.

## B.2 Human Agreement with Automatic Judgment

To further validate the reliability of LLM-based factual evaluation, we conduct a manual inspection on a randomly sampled subset of 100 responses generated by **MemoryOS**. Each response is independently annotated by human evaluators for factual correctness and compared against the binary judgments produced by DeepSeek-V3.2. We focus on MemoryOS as it represents the strongest memory-augmented agent and therefore constitutes the most critical case for validating evaluation reliability.

Figure 4 shows the resulting confusion matrix between human annotations and DeepSeek-based judgments. Out of 100 samples, 86 cases receive consistent labels, corresponding to an agreement rate of 86%. Disagreements primarily arise in borderline cases involving partial recall or implicit temporal assumptions, where factual validity depends on nuanced interpretation rather than explicit contradiction. Overall, this analysis supports the use of LLM-as-a-Judge as a reliable proxy for large-scale comparative evaluation in the NaturalMem setting.

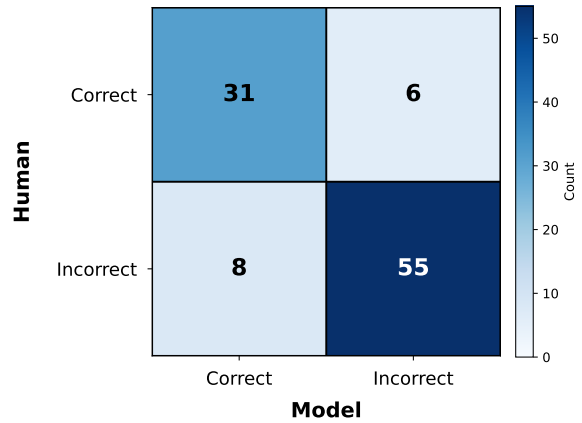


Figure 4: Confusion matrix comparing human annotations and DeepSeek-V3.2 factual judgments on 100 randomly sampled responses generated by MemoryOS. Diagonal entries indicate agreement between human evaluators and the LLM-based judge, while off-diagonal entries correspond to disagreement cases.

## C Prompt Design for Evaluation and Agent Responses

This appendix documents the prompt templates used in our experiments for (1) factual correctness evaluation, (2) style consistency evaluation, and (3) response generation by memory-augmented agents. These prompts are designed to be minimal, explicit, and task-specific, in order to reduce ambiguity and avoid unintended instruction leakage.

For evaluation, we adopt a binary factual judgment protocol and a graded stylistic alignment protocol, following prior LLM-as-a-Judge practices. For memory-augmented agents, prompts are restricted to providing retrieved memory and the current user query, without additional task-specific instructions.

Someone is talking with you.  
 You have these conversation memories: {context}  
 At {query\_timestamp}, they talk with you about: {question}  
 How would you respond?

Figure 5: Prompt used by A-Mem to generate responses, where retrieved memory is provided together with the current user query.

Your task is to label an answer as CORRECT or WRONG based on whether it satisfies the constraint.

**You will be given:**

1. The question
2. The constraint for response
3. The generated answer

**Evaluation rules:**

- Be generous: if the answer clearly respects the constraint, mark CORRECT.
- If the answer violates any explicit restriction, mark WRONG.
- Style does NOT matter, only factual and logical compliance. Please provide an objective and fair evaluation.

**Now evaluate:**

Question: {question}

Constraint: {constraint}

Generated Answer: {generated\_answer}

First, provide a one-sentence explanation.

Then output CORRECT or WRONG and reason.

**Return JSON only:**

```
{{  
  "label": "CORRECT or WRONG"  
  "reason": "<short explanation>"  
}}
```

Output MUST be exactly one JSON object and nothing else.

Figure 6: Prompt used for factual correctness evaluation, where the judge determines whether a generated answer satisfies an explicit constraint, ignoring stylistic quality.

You are an expert evaluator of conversational style.

**You will be given:**

- A user question
- An assistant answer
- A source\_show and an assistant character

Your task is to evaluate whether the assistant's answer matches the given assistant character.

**Important note:** The assistant character may come from well-known television shows, movies, or other media (e.g., characters from Friends, The Big Bang Theory, or similar series). You are allowed and expected to use your prior knowledge about these characters (their personality traits, typical speaking style, emotional patterns, and behavior) to make an informed judgment.

**You MUST evaluate from the following dimensions:**

1. Tone and attitude (polite, warm, professional, casual, etc.)
2. Speaking style (verbosity, formality, emotional expression)
3. Role consistency (does it behave like the described assistant?)
4. Overall stylistic alignment Please provide an objective and fair evaluation.

**Give a score from 1 to 5:**

- 5 = Perfectly matches the assistant character
- 4 = Mostly matches with minor inconsistencies
- 3 = Partially matches
- 2 = Poor match
- 1 = Completely inconsistent

Question: {question}

source\_show: {source\_show}

assistant\_character: {assistant\_character}

Answer: {answer}

**Only output a JSON object in the following format:**

```
{{
  "score": <integer from 1 to 5>
  "reason": "<short explanation>"
}}
```

Figure 7: Prompt used for style consistency evaluation, assessing alignment with a target assistant character on a 1–5 scale.