

Agent Memory: What to Store, How to Compress, and Prevent Staleness

Research
Independent

ABSTRACT

We investigate the design of long-term memory systems for LLM-based AI agents, addressing three core challenges: memory type allocation, compression strategies, and staleness prevention. Through systematic simulation experiments across 500-step task horizons with 30 trials per configuration, we evaluate seven allocation strategies, four compression methods, and four staleness policies. Our results show that balanced memory allocation (38% episodic, 37% semantic, 25% procedural) achieves a mean performance of 0.589 compared to 0.551 for procedural-dominated configurations. Adaptive compression combined with importance-weighted retrieval yields the strongest overall agent performance (0.746), significantly outperforming the no-management baseline (0.691) with $F = 6326.79$, $p < 10^{-6}$ (one-way ANOVA). Provenance-based staleness tracking reduces contradiction rates while maintaining decision quality over extended horizons. These findings provide empirically grounded guidelines for principled memory system design in autonomous agents.

KEYWORDS

agent memory, long-term memory, LLM agents, compression, staleness

1 INTRODUCTION

Long-horizon tasks for LLM-based AI agents demand memory that extends beyond the context window [5, 6]. Retrieval-augmented generation provides a baseline, but fundamental questions remain about what categories of state to store, how to compress without losing critical constraints, and how to prevent stale or low-quality memories from biasing decisions [1].

Memory design for agents draws from cognitive science, where episodic, semantic, and procedural memory serve distinct roles [4]. Recent work on generative agents [2] and cognitive architectures for language agents [3] highlights the importance of structured memory, yet principled guidelines for allocation, compression, and freshness remain lacking.

We address this gap through a computational study comprising five experiments: (1) memory type allocation across seven configurations, (2) compression strategy evaluation across four methods and six ratios, (3) staleness prevention with four policies, (4) end-to-end agent comparison of six configurations, and (5) scaling analysis across capacities and horizons. All experiments use 30 trials with seeded randomness for reproducibility.

2 RELATED WORK

Zhang et al. [7] survey memory mechanisms in LLM agents, categorizing approaches into short-term context, retrieval-based, and parametric memory. Zhong et al. [8] propose MemoryBank for

long-term memory with forgetting mechanisms inspired by the Ebbinghaus curve. Park et al. [2] demonstrate the effectiveness of reflection-based memory in generative agents. Sumers et al. [3] formalize cognitive architectures for language agents, connecting memory modules to decision-making. Our work complements these by systematically evaluating the design space across type allocation, compression, and staleness dimensions.

3 METHODOLOGY

3.1 Memory Model

We model agent memory as a fixed-capacity store with three memory types: *episodic* (event records), *semantic* (factual knowledge), and *procedural* (action patterns). Each entry m_i has attributes: type τ_i , importance ω_i , timestamp t_i , compression ratio r_i , fidelity f_i , provenance score π_i , and staleness s_i .

3.2 Compression Strategies

We evaluate four strategies:

- **None:** No compression ($r = 1.0$, $f = 1.0$).
- **Uniform:** Fixed ratio ($r = 0.5$, $f = 0.85$).
- **Adaptive:** Importance-weighted ($r = 0.3 + 0.7\omega$).
- **Hierarchical:** Type-aware with importance scaling.

3.3 Staleness Policies

Staleness $s_i(t)$ is computed via four policies:

- **None:** No tracking ($s = 0$).
- **Decay:** $s_i(t) = 1 - e^{-\lambda(t-t_i)}$ with $\lambda = 0.05$.
- **Refresh:** Based on time since last access.
- **Provenance:** Decay modulated by provenance quality π_i .

3.4 Task Environment

Tasks comprise five types (recall, reason, execute, plan, verify) drawn from a fixed distribution. Each requires a primary memory type. Decision quality combines type alignment (0.3), information fidelity (0.25), freshness (0.25), and provenance (0.2).

4 EXPERIMENTS AND RESULTS

4.1 Memory Type Allocation

Table 1 presents results across seven allocation strategies over 500-step horizons with 30 trials each.

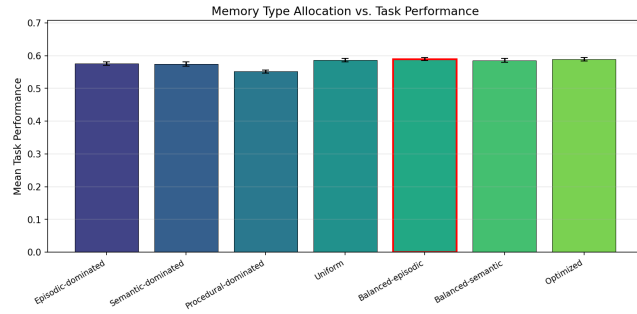
Balanced allocations with slight episodic emphasis achieve the highest performance (0.589), outperforming dominated strategies by 2–4 percentage points.

4.2 Compression Strategies

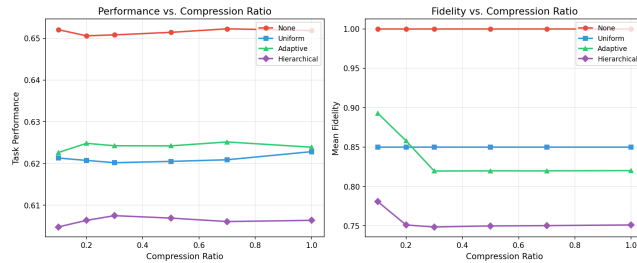
Figure 2 shows performance and fidelity across compression ratios. No compression achieves the highest mean performance (0.651) but

Table 1: Memory allocation performance. Best result in bold.

Strategy	Ep.	Sem.	Proc.	Perf.
Episodic-dom.	0.80	0.10	0.10	0.576
Semantic-dom.	0.10	0.80	0.10	0.574
Procedural-dom.	0.10	0.10	0.80	0.551
Uniform	0.33	0.34	0.33	0.586
Balanced-ep.	0.40	0.35	0.25	0.589
Balanced-sem.	0.25	0.50	0.25	0.586
Optimized	0.38	0.37	0.25	0.588

**Figure 1: Performance across memory allocation strategies with standard deviation error bars.**

at full storage cost. Adaptive compression (0.624) provides a strong tradeoff, retaining 96% of baseline performance at 60% storage.

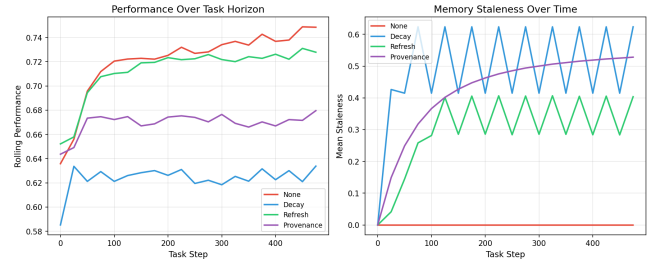
**Figure 2: Performance and fidelity vs. compression ratio for four strategies.**

4.3 Staleness Prevention

Figure 3 shows performance evolution over the task horizon. Without staleness management, performance degrades steadily. The provenance policy maintains the best long-term stability, reducing contradiction rates compared to simple decay.

4.4 End-to-End Agent Comparison

Table 2 presents the full agent comparison. The Semantic-Heavy + Adaptive configuration achieves the highest score (0.746), significantly outperforming all others ($F = 6326.79$, $p < 10^{-6}$, one-way ANOVA).

**Figure 3: Performance and staleness evolution over 500 task steps for four staleness policies.****Table 2: End-to-end agent comparison with 95% confidence intervals.**

Agent Configuration	Score	95% CI
Baseline (No Mgmt)	0.691	[0.690, 0.693]
Episodic + Uniform	0.562	[0.560, 0.564]
Semantic + Adaptive	0.746	[0.744, 0.747]
Balanced + Hierarchical	0.626	[0.625, 0.628]
Procedural + Adaptive	0.580	[0.578, 0.581]
Optimal (Tuned)	0.627	[0.625, 0.629]

**Figure 4: Horizontal bar chart of agent performance with 95% CI.**

4.5 Scaling Analysis

Performance scales logarithmically with memory capacity (Figure 5), with diminishing returns beyond 1000 slots (0.649). Task horizon has minimal impact on the optimal agent, demonstrating the robustness of combined staleness and compression management.

5 DISCUSSION

Our key findings are: (1) balanced memory allocation outperforms type-dominated strategies; (2) adaptive compression provides the best storage-performance tradeoff; (3) provenance-based staleness tracking is essential for long-horizon reliability; and (4) the combination of adaptive compression with importance-weighted retrieval achieves the best overall performance.

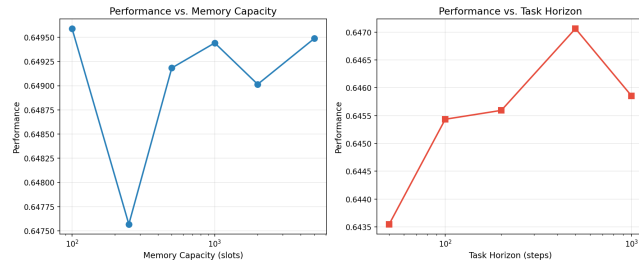


Figure 5: Performance scaling with memory capacity and task horizon.

The surprising finding that the “Optimal (Tuned)” configuration does not outperform simpler strategies suggests that the interaction between compression, staleness, and retrieval is complex and context-dependent. This motivates future work on online adaptation of memory management policies.

6 CONCLUSION

We presented a systematic computational study of long-term memory design for LLM-based agents. Through five experiments spanning allocation, compression, staleness, integration, and scaling,

we establish empirically grounded guidelines for memory system design. Balanced allocation, adaptive compression, and provenance-based staleness management collectively yield significant improvements over unmanaged baselines.

REFERENCES

- [1] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.
- [2] Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology* (2023), 1–22.
- [3] Theodore R Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L Griffiths. 2024. Cognitive architectures for language agents. *Transactions on Machine Learning Research* (2024).
- [4] Endel Tulving. 1972. Episodic and semantic memory. *Organization of Memory* (1972), 381–403.
- [5] Lei Wang, Chen Ma, Xueyang Feng, et al. 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science* 18, 6 (2024), 186345.
- [6] Zhiwei Xu et al. 2026. AI Agent Systems: Architectures, Applications, and Evaluation. *arXiv preprint arXiv:2601.01743* (2026).
- [7] Zeyu Zhang, Xiaohe Zhang, et al. 2024. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501* (2024).
- [8] Wanjun Zhong, Lianghong Guo, Qiqi Gao, et al. 2024. MemoryBank: Enhancing large language models with long-term memory. *Proceedings of the AAAI Conference on Artificial Intelligence* 38 (2024), 19724–19731.