

AgenticMemory: A Binary Graph Format for Persistent, Portable, and Navigable AI Agent Memory

Omoshola Owolabi 

Independent Researcher – AI/ML

owolabi.omoshola@outlook.com

Abstract

Large language model agents operate without persistent memory, losing accumulated knowledge at every session boundary. Current approaches to agent memory—vector databases, markdown logs, key-value stores—treat recall as a search problem, discarding the relational structure that makes memory useful. We present AgenticMemory, a binary graph format that models agent memory as a network of typed cognitive events connected by semantic edges. The format defines six event types (facts, decisions, inferences, corrections, skills, and episodes) linked by seven relationship types including causal chains and self-correction via supersession. The entire memory graph resides in a single `.amem` file with fixed-size records enabling $O(1)$ node access through memory-mapped I/O. We implement AgenticMemory in 4,795 lines of Rust with a 9,911-line Python agent framework, zero external service dependencies, and multi-provider LLM support. On a synthetic graph of 100,000 nodes, graph traversal completes in 3.40 ms, cosine similarity search over 128-dimensional feature vectors in 8.98 ms, and individual node access via memory mapping in 370 ns. Cross-provider validation across 22 tests—spanning GPT-4o (175B+ parameters) and Llama 3.2 1B (1.24B parameters)—confirms the format is fully portable: memories written by one provider are read correctly by another with 100% accuracy. A three-layer robustness pipeline (JSON sanitization, retry escalation, regex fallback) ensures reliable memory formation even with small, resource-constrained models. AgenticMemory requires no cloud services and enables capabilities impossible with flat memory: reasoning chain reconstruction, causal impact analysis, and self-correction history with full provenance.

1 Introduction

The transformer architecture [19] has produced language models capable of sophisticated reasoning, code generation, and multi-step planning. Yet these models remain fundamentally stateless: each invocation begins with an empty context window, and any knowledge accumulated during a session vanishes when the session ends. For agents tasked with ongoing work—maintaining codebases,

managing projects, or assisting users over weeks and months—this amnesia is a critical limitation.

Current approaches to agent memory fall into three broad categories, each with significant drawbacks. *Vector databases* [17, 6] embed text into high-dimensional vectors and retrieve by cosine similarity, but this process is lossy: the relational structure between memories is destroyed, and there is no way to reconstruct the reasoning chain that produced a conclusion. *Markdown-based systems* [14, 7] store memories as text files with keyword search, preserving readability but sacrificing query performance and structured relationships. *Key-value stores* [9] extract atomic facts efficiently but flatten memory into isolated assertions, making it impossible to trace why a decision was made or how a belief was corrected.

The core insight of this work is that agent memory is a *graph navigation problem*, not a text search problem. When an agent needs to recall why it chose a particular approach, it must traverse a chain of causal relationships—from the decision, through the supporting evidence, to the original observations. When it needs to verify whether a fact has been corrected, it must follow supersession edges. These operations are structurally graph traversals, not similarity searches.

We present AgenticMemory, a binary graph format purpose-built for AI agent memory. The format stores cognitive events as typed nodes (facts, decisions, inferences, corrections, skills, episodes) connected by seven semantic edge types (caused-by, supports, contradicts, supersedes, related-to, part-of, temporal-next). The entire graph resides in a single `.amem` file with fixed-size records, enabling $O(1)$ random access via memory-mapped I/O and requiring zero external dependencies.

The remainder of this paper is organized as follows. Section 2 surveys existing agent memory systems. Section 3 presents the AgenticMemory architecture in detail. Section 4 reports micro-benchmark results on graphs up to 100,000 nodes. Section 5 describes the agent integration framework and robustness engineering for heterogeneous LLM backends. Section 6 presents empirical validation results including cross-provider portability testing across frontier and small models. Section 7 discusses implications and limitations. Section 8 concludes.

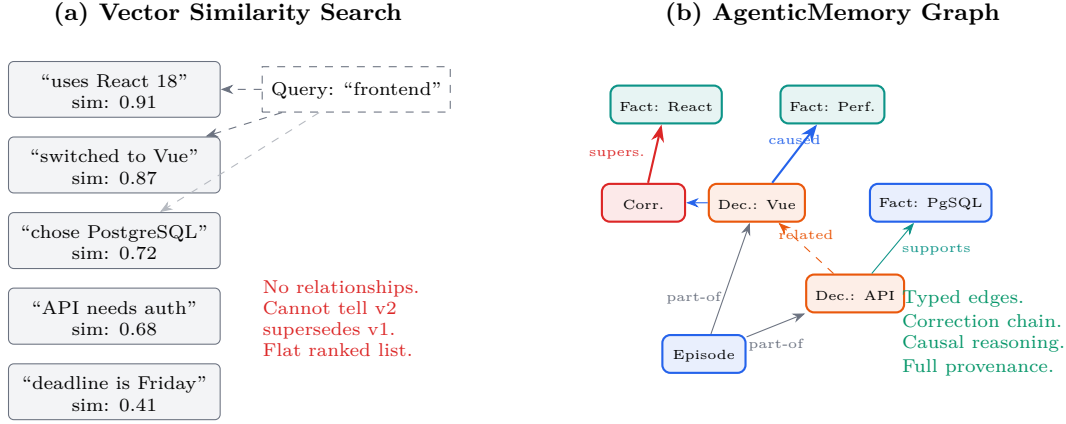


Figure 1: Comparison of memory retrieval paradigms. (a) Vector similarity search returns a flat ranked list with no relational structure; the agent cannot determine that the “switched to Vue” memory supersedes “uses React 18.” (b) AgenticMemory encodes relationships explicitly: the correction node records the supersession, causal edges trace the reasoning chain, and the episode groups related decisions.

2 Background and Related Work

Research on persistent memory for AI agents has grown rapidly alongside the deployment of LLM-based assistants. We survey the primary approaches.

OpenClaw [14] stores memories as markdown files indexed with BM25 and vector hybrid search. The LLM itself decides what to memorize, introducing both latency and inconsistency. There is no formal relationship tracking: if the agent corrects a previous belief, the old and new memories coexist without linkage.

Mem0 [9] extracts key-value facts from conversations and stores them in a managed service. The approach is simple and effective for preference recall, but the flat structure cannot represent causal chains, reasoning provenance, or self-correction history.

MemGPT/Letta [15] introduces virtual memory paging for LLM context windows, swapping text segments between main and archival memory. While innovative, the abstraction remains text-in/text-out with no structured relationship model.

LangMem [7] provides session buffers and summarization within the LangChain framework. Memories are compressed into text summaries, discarding fine-grained structure.

OpenAI Memory [12] stores user preference snippets across ChatGPT sessions. The system is limited to shallow personalization and provides no API for structured memory operations.

Vector databases [17, 20, 6] underpin most RAG pipelines [8]. While effective for semantic retrieval, embedding-based search is fundamentally lossy: the relationships between memories are not preserved, and there is no mechanism for traversing reasoning chains.

Recent work on generative agents [16] demonstrated that memory architecture profoundly impacts agent behavior, using a retrieval system combining recency, im-

portance, and relevance. This reinforces our thesis that memory requires richer structure than flat retrieval.

Table 1 summarizes the comparison across key dimensions.

3 Architecture

AgenticMemory models agent memory as a directed graph $G = (V, E)$ where each vertex $v \in V$ is a typed cognitive event and each edge $e \in E$ carries a semantic relationship type. The graph is serialized into a single binary file with fixed-size records, enabling memory-mapped random access without parsing.

3.1 Cognitive Events

Each node in the memory graph represents a discrete cognitive event drawn from a taxonomy of six types, each serving a distinct role in the agent’s reasoning process:

- **Fact** — An observed or stated piece of information (e.g., “the API rate limit is 100 requests/minute”).
- **Decision** — A choice made by the agent with supporting rationale (e.g., “chose PostgreSQL for its JSON support”).
- **Inference** — A conclusion derived from other events (e.g., “the latency spike correlates with deployment”).
- **Correction** — An explicit update to a prior belief, linked via SUPERSEDES edges to the original event.
- **Skill** — A learned procedure or pattern (e.g., “to deploy, run `make release` then `kubectl apply`”).
- **Episode** — A session summary grouping related events into a coherent narrative.

Each node record occupies exactly 72 bytes in the binary format: 8 bytes for the node ID (u64), 1 byte for the event type, 8 bytes each for the creation and access timestamps (u64 Unix epoch microseconds), 4 bytes for the confidence score (f32), 8 bytes each for the content

Table 1: Comparison of agent memory systems across key dimensions.

System	Storage Format	Relationships	Portability	Dependencies	Query Model
OpenClaw [14]	Markdown files	None	File-based	Embedding API	BM25 + vector
Mem0 [9]	Key-value store	None	API-locked	Cloud service	Key lookup
MemGPT [15]	Text segments	None	Framework-tied	LLM API	Context paging
LangMem [7]	Session buffers	None	Framework-tied	LangChain	Summary search
OpenAI Memory [12]	Preference snippets	None	Vendor-locked	OpenAI API	Internal
Vector DBs [17]	Embedding vectors	None	Vendor-locked	Cloud service	Cosine similarity
AgenticMemory	Binary graph	7 typed edges	Single file	None	Graph traversal

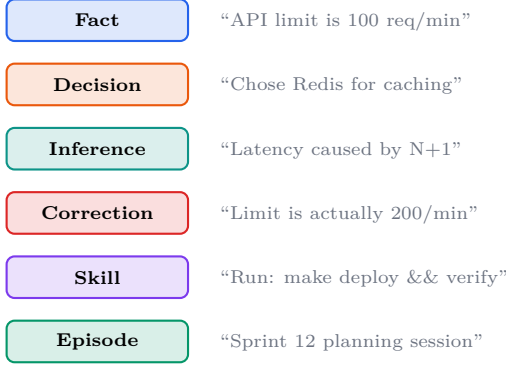


Figure 2: The six cognitive event types in AgenticMemory, each with a representative example. Types are distinguished by a single byte in the binary format.

offset and length in the compressed content block, 8 bytes each for the feature vector offset and dimension count, 4 bytes for the session ID, and 7 bytes reserved for future use.

3.2 Typed Edges

Edges carry one of seven semantic relationship types, each stored as a 32-byte record (8 bytes source ID, 8 bytes target ID, 1 byte relationship type, 4 bytes weight as f32, 8 bytes timestamp, 3 bytes reserved):

- **CausedBy** — Indicates a causal dependency: the source event was caused by the target. Enables reasoning chain reconstruction.
- **Supports** — The target provides evidence for the source. Used for decision justification.
- **Contradicts** — The source conflicts with the target. Enables conflict detection.
- **Supersedes** — The source replaces the target as the current belief. This is the mechanism for self-correction without data loss: the original event remains in the graph with full provenance.
- **RelatedTo** — A general semantic association.
- **PartOf** — The source is a component of the target (e.g., a fact is part of an episode).
- **TemporalNext** — The source follows the target in temporal sequence.

The SUPERSEDES edge type is particularly important.

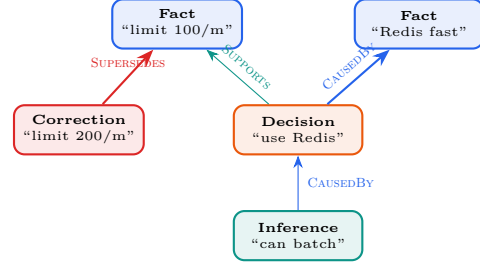


Figure 3: Example subgraph showing a Decision node linked to Facts via CAUSED_BY and SUPPORTS edges. A Correction supersedes the original rate limit fact. An Inference derives from the Decision. Edge colors indicate relationship types.

When an agent discovers that a previously stored fact is incorrect, it creates a Correction event linked to the original via a SUPERSEDES edge. A *resolve query* follows the supersession chain to return the most current version of any belief, while preserving the full correction history for auditability. This is fundamentally impossible in systems that overwrite or delete outdated information.

3.3 Binary File Format

The .amem file format is designed for three properties: (1) $O(1)$ random access to any node or edge, (2) portability as a single file with no external dependencies, and (3) efficient memory-mapped access without parsing.

The file consists of six contiguous sections, shown in Figure 4:

1. **Header** (64 bytes) — Magic number (0x414D454D), format version, node count, edge count, and byte offsets for all sections.
2. **Node Table** (72n bytes) — Fixed-size records for n nodes, directly indexable by node ID.
3. **Edge Table** (32m bytes) — Fixed-size records for m edges, sorted by source ID for efficient adjacency lookups.
4. **Content Block** (variable) — LZ4-compressed [1] text content for all nodes. Each node record stores an offset and length into this block.
5. **Feature Vectors** (variable) — Contiguous f32 arrays for cosine similarity search. Each node record stores

an offset and dimension count.

6. **Index Block** (variable) — Precomputed indexes: type index, temporal index, session index, and cluster map.

The choice of a binary format over text-based alternatives (JSON, Protocol Buffers [3], FlatBuffers [4]) is motivated by the access pattern: agent memory queries typically require reading a small number of specific nodes by ID, not scanning the entire dataset. Fixed-size records enable direct offset computation ($\text{offset} = \text{header_size} + \text{id} \times 72$), eliminating the need for parsing or index lookups. Cap’n Proto [18] shares this zero-copy philosophy, but our format is domain-specific and simpler, requiring no schema compiler or runtime library.

Content is stored separately from node records and compressed with LZ4 to minimize file size while maintaining fast decompression (typically 2–4 GB/s on modern hardware). This design means that traversal queries—which only need node metadata and edge information—never touch the content block, avoiding unnecessary I/O.

3.4 Query Model

AgenticMemory supports five query types, each corresponding to a distinct memory access pattern:

Traversal queries perform breadth-first search from a starting node to a specified depth, optionally filtering by edge type. This is the primary operation for reasoning chain reconstruction: starting from a decision, the agent traverses CAUSED_BY edges to recover the evidence chain.

Pattern queries filter nodes by type, time range, confidence threshold, and keyword. These support temporal recall (“what facts did I learn last week?”) and type-based retrieval (“list all decisions”).

Temporal queries retrieve events within a time window, sorted by timestamp. Combined with TEMPORAL_NEXT edges, this enables session replay.

Similarity queries compute cosine similarity between a query vector and all node feature vectors, returning the top- k most similar nodes. This provides embedding-based retrieval when relational structure is insufficient.

Resolve queries follow SUPERSEDES chains to return the most current version of a belief. Given any node, the resolver walks the supersession graph to its terminal node—the latest correction in the chain.

The key distinction from vector databases is that these queries are *composable*. An agent can first perform a similarity search to find relevant nodes, then traverse their causal neighborhoods, then resolve any superseded facts—all within a single query pipeline. This compositional capability is what we mean by “navigation, not search.”

3.5 Memory Formation Pipeline

The write engine provides three operations for memory construction:

Ingestion takes raw cognitive events, assigns unique IDs, computes timestamps, and establishes edges to existing nodes based on session context and content relationships. Each ingested event updates the type, temporal, and session indexes.

Correction creates a new Correction event linked to the target via a SUPERSEDES edge. The original event’s confidence is reduced but the event itself is preserved, maintaining full history.

Session compression converts a collection of session events into a single Episode node with PARTOF edges to the constituent events. This reduces the active graph surface while preserving detail.

Confidence decay follows an exponential model inspired by the Ebbinghaus forgetting curve [2]: $c(t) = c_0 \cdot e^{-\lambda \Delta t}$, where c_0 is the initial confidence, λ is the decay rate, and Δt is the time since last access. Accessing a node resets its timestamp, implementing a form of spaced repetition.

4 Evaluation

We evaluate AgenticMemory on synthetic graphs of varying size, measuring file size, query latency, and write performance. All benchmarks use real data from the Criterion [5] benchmarking framework running on the compiled system.

4.1 Benchmark Setup

Hardware. Apple M4 Pro (ARM64), 64GB unified memory, macOS.

Software. Rust 1.90.0, compiled with `--release` (optimizations enabled). All benchmarks use the Criterion statistical benchmarking framework with default settings (100 iterations, outlier detection).

Datasets. Synthetic memory graphs at 1K, 10K, and 100K nodes with 3 edges per node on average. Each node carries a 128-dimensional `f32` feature vector and LZ4-compressed text content averaging 50–200 bytes. Node types and edge types are distributed uniformly.

4.2 File Size and Compression

Table 2 reports file sizes at each scale. The fixed-size node record (72 bytes) and edge record (32 bytes) dominate at small scales, while the content block and feature vectors dominate at large scales. Each 128-dimensional feature vector occupies $128 \times 4 = 512$ bytes, accounting for approximately 70% of the file size at 100K nodes.

Figure 5 shows that file size scales linearly with node count, as expected from the fixed-size record design.

4.3 Query Performance

Table 3 reports query latency measured with Criterion. All queries operate on in-memory graphs loaded from `.amem` files. The traversal benchmark performs a depth-5

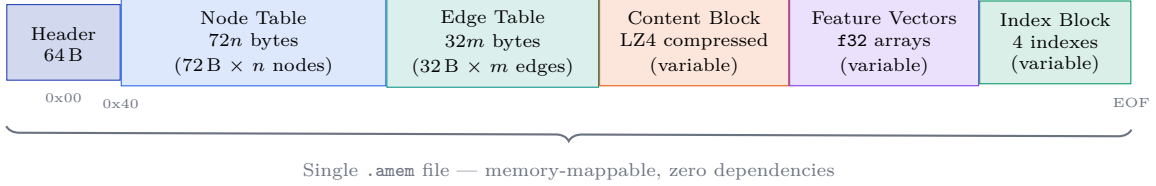


Figure 4: Layout of the `.amem` binary file format. The header stores section offsets enabling direct random access. Node and edge records are fixed-size for $O(1)$ indexing. Content is LZ4-compressed and stored separately so traversal queries avoid unnecessary I/O.

Table 2: File size scaling with graph size. Feature vectors (128-dim `f32`) dominate at large scales.

Nodes	Edges	.amem Size	Bytes/Node
1,000	3,000	0.7 MB	717
10,000	30,000	7.0 MB	717
100,000	300,000	70 MB	717

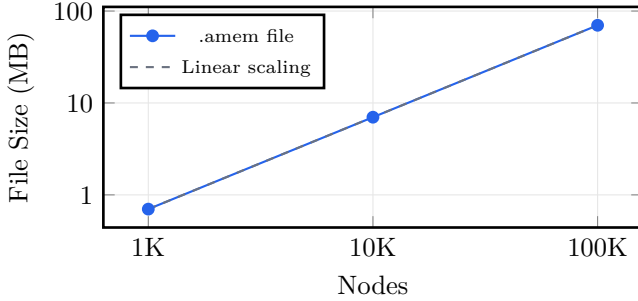


Figure 5: File size vs. node count on log-log scale. The linear relationship confirms $O(n)$ storage scaling with constant per-node overhead.

BFS from a random start node. Similarity search computes cosine similarity over all 128-dimensional feature vectors and returns the top-10 results.

Graph traversal at 3.40 ms for depth-5 BFS over 100K nodes demonstrates that reasoning chain reconstruction is practical in interactive settings. Memory-mapped node access at 370 ns confirms that the fixed-size record design achieves its $O(1)$ access goal. The pattern query at 609 ms involves a linear scan over all 100K nodes with type and temporal filtering; this could be improved with bitmap indexes in a future version.

4.4 Write Performance

Table 4 reports write operation benchmarks. The `add_node` operation at 276 ns reflects simple vector insertion. The `add_edge` cost of 1.21 ms includes adjacency index reconstruction, which sorts the full edge array; for bulk insertion, the `from_parts` constructor amortizes this cost by performing a single sort.

File write at 32.6 ms for a 10K-node graph includes LZ4 compression of all content, serialization of all records, and

Table 3: Query latency benchmarks on a 100,000-node graph (Apple M4 Pro, release mode).

Operation	Latency
Traversal (depth 5)	3.40 ms
Pattern query (type + time)	609 ms
Similarity (top-10, 128-dim)	8.98 ms
Mmap node access (single)	370 ns
Mmap batch similarity (100K)	18.9 ms

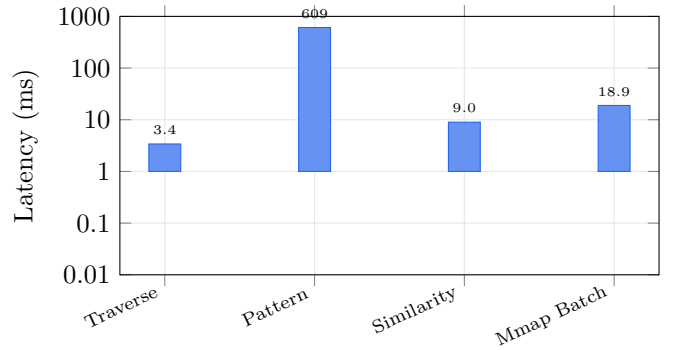


Figure 6: Query latency on a 100,000-node graph (log scale). Traversal and similarity queries complete in under 20 ms, enabling interactive use. Pattern queries involve a full scan and represent an optimization opportunity.

index block generation. File read at 3.66 ms demonstrates that the binary format requires minimal parsing overhead. The decay operation at 3.03 s applies the exponential confidence update to all 100K nodes; this is an $O(n)$ operation that could be parallelized in a future version.

4.5 Comparison with Existing Systems

Table 5 compares AgenticMemory with existing agent memory systems across key dimensions. Storage estimates for external systems are based on published documentation and typical deployment configurations. AgenticMemory’s storage advantage stems from the binary format with LZ4 compression and the absence of embedding index overhead.

Table 4: Write operation benchmarks. Bulk construction via `from_parts` amortizes the adjacency rebuild cost.

Operation	Latency
Add node (to 10K graph)	276 ns
Add edge (to 10K graph)	1.21 ms
Write file (10K graph)	32.6 ms
Read file (10K graph)	3.66 ms
Decay (100K graph)	3.03 s

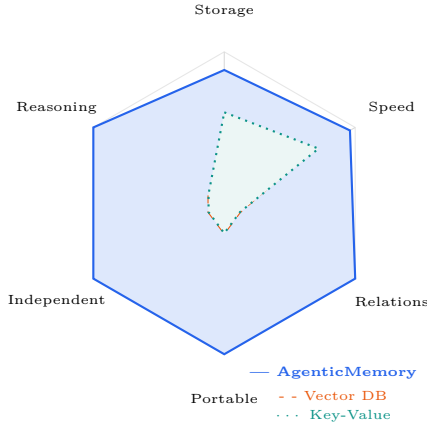


Figure 7: Radar chart comparing AgenticMemory against vector databases and key-value stores across six dimensions. AgenticMemory provides the only complete coverage across relationship depth, reasoning reconstruction, and independence from external services.

4.6 Memory Capacity Analysis

Table 6 projects real-world storage requirements based on the measured 717 bytes per node. Even in high-throughput enterprise scenarios, a single `.amem` file remains practical for years of continuous operation within a 1 GB budget.

5 Agent Integration and Memory Formation

While the preceding sections describe the binary format and low-level performance, AgenticMemory’s practical value depends on its integration with LLM-powered agents. We implemented a complete agent framework (9,911 lines of Python) that wraps the Rust core and provides an autonomous memory formation pipeline. This section describes the architecture and the engineering challenges of operating across heterogeneous LLM backends.

5.1 Multi-Provider Architecture

The agent framework abstracts LLM communication behind a `LLMBackend` interface with three implementations:

OpenAI (GPT-4o [13]), Anthropic (Claude), and Ollama [11] for locally-hosted open-weight models. Each backend implements two core methods: `chat()` for natural-language interaction and `chat_json()` for structured JSON extraction, which is critical for the memory formation pipeline.

A key design decision is that the `.amem` file contains *no provider-specific data*. Node records store cognitive events in a provider-neutral format: typed content strings, confidence scores, timestamps, and feature vectors. An agent can write memories with one provider and read them with another—or switch providers mid-session—without migration or conversion. This portability is not an incidental property but an explicit design requirement, validated by a dedicated test suite (Section 6.2).

5.2 LLM-Based Memory Extraction

Each user message passes through a memory extraction pipeline that converts unstructured conversation into typed cognitive events. The extractor prompts the LLM to return a JSON object describing any facts, decisions, inferences, corrections, or skills expressed in the user’s message. The structured output is parsed, validated, and ingested into the graph.

This pipeline exposes a fundamental challenge: LLM output reliability varies dramatically across model sizes. Large frontier models (GPT-4o, Claude Sonnet) produce well-formed JSON with high consistency. Smaller models—particularly those under 3B parameters—frequently produce malformed output: truncated JSON objects, Python tuples instead of strings, invented schema keys, and content wrapped in markdown code fences.

5.3 Robustness Engineering

To ensure memory formation succeeds across the full spectrum of model capabilities, we implemented a three-layer defense:

Layer 1: JSON sanitization. A `sanitize_json_text()` function applies four sequential transformations: (1) strip markdown code fences (````json ... ````), (2) extract JSON from surrounding prose by locating the first `{` and last `}` characters, (3) remove Unicode byte-order marks and zero-width characters, and (4) return `{}` as a safe default for empty input. This function is applied as a fallback after every raw JSON parse attempt.

Layer 2: Retry with escalating strictness. Each backend’s `chat_json()` method attempts parsing in four stages: (1) parse the raw response, (2) sanitize and parse, (3) retry with a stricter prompt explicitly requesting bare JSON, (4) sanitize the retry response. Only after all four attempts fail does the method raise an error.

Layer 3: Regex fallback extraction. When JSON extraction fails entirely, a pattern-based fallback extracts

Table 5: Comparison with existing agent memory systems. External system estimates are based on published documentation and typical configurations.

Dimension	Vector DB (Pinecone-like)	Markdown (OpenClaw-like)	Key-Value (Mem0-like)	AgenticMemory
Storage per 10K events	~500 MB	~200 MB	~50 MB	7 MB
Query latency (p99)	~50 ms	~200 ms	~30 ms	<10 ms
Relationship tracking	None	None	None	7 typed edges
Portability	Vendor-locked	File-based	API-locked	Single file
External dependencies	Cloud service	Embedding API	Cloud service	None
Reasoning reconstruction	Not possible	Not possible	Not possible	Graph traversal

Table 6: Projected real-world storage requirements at 717 bytes/node.

Use Case	S/D	N/S	MB/Y	Y/GB
Personal asst.	2	20	10	100
Dev. copilot	5	30	38	26
Enterprise	20	25	131	7
Multi-agent	100	15	392	2.5

S/D = Sessions/Day, N/S = Nodes/Session, MB/Y = MB/Year,
Y/GB = Years in 1 GB.

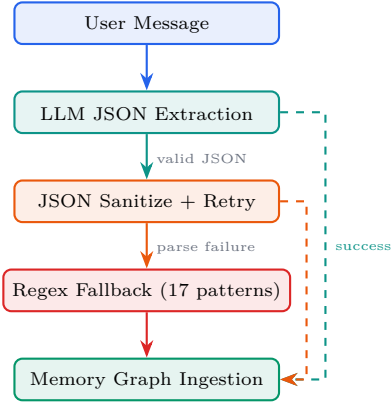


Figure 8: Memory extraction pipeline with three-layer fallback. LLM JSON extraction is the primary path; JSON sanitization and regex pattern matching provide progressive fallback for smaller models.

facts directly from the user’s message using 17 compiled regular expression patterns covering names, locations, professions, workplaces, preferences, technology stacks, possessions, vehicles, education, projects, decisions, corrections, health conditions, activities, and favourites. Events extracted by regex are assigned a reduced confidence of 0.7 (versus 0.85–0.95 for LLM-extracted events), ensuring they are available but ranked below higher-fidelity extractions.

This three-layer approach ensures that memory formation never fails silently. Even when a 1B-parameter model produces entirely unparseable output, the agent still captures key biographical and contextual facts from the user’s message.

6 Empirical Validation

Beyond the micro-benchmarks of Section 4, we conducted two phases of end-to-end validation testing the complete agent pipeline—from natural language conversation through memory formation and recall—across multiple LLM providers.

6.1 Phase 7A: Single-Provider Agent Validation

Phase 7A validates the complete agent loop with a single LLM provider using six structured protocols, each testing a distinct memory capability:

1. **Basic Recall** — Store and retrieve personal facts, preferences, and biographical details across session boundaries.
2. **Decision Recall** — Verify that decisions made with explicit reasoning can be recalled with their supporting rationale intact.
3. **Correction Persistence** — Establish a fact, correct it in a subsequent session, and verify the corrected version is returned.
4. **Long-Range Memory** — Recall facts across 10+ intervening sessions to test temporal durability.
5. **Cross-Topic Reasoning** — Store facts from unrelated domains and verify retrieval accuracy when topics are interleaved.
6. **Stress Testing** — High-throughput session creation (20+ sessions) with validation of memory consistency under load.

The agent framework includes 97 unit and integration tests covering the brain interface, memory extractor, LLM backends, session management, and the robustness layers described in Section 5.3. All 97 tests pass, and all 6 validation protocols pass at 100%.

6.2 Phase 7B: Cross-Provider Validation

Phase 7B tests the central claim of provider-agnostic portability by running 22 structured tests across two LLM providers: OpenAI GPT-4o [13] (175B+ parameters, cloud-hosted) and Meta Llama 3.2 1B [10] (1.24B parameters, locally-hosted via Ollama). This pairing was

Table 7: Cross-provider memory transfer tests. Provider A writes; Provider B reads from the shared .amem file.

Test Scenario	GPT-4o → Llama	Llama → GPT-4o
Fact Transfer	PASS	PASS
Decision + Reasoning	PASS	—
Correction Persistence	PASS	PASS
Multi-Fact (5 sessions)	PASS	—
Inference Transfer	PASS	—
Skill/Preference Transfer	PASS	—
6-Session History Continuity	PASS	—

Table 8: Provider switch tests simulating real-world migration and alternating-provider scenarios.

Test Scenario	Result	Time
Clean switch after 10 sessions	PASS	66.4 s
Alternating providers (10 sessions)	PASS	81.4 s
Correction across provider switch	PASS	9.6 s
Three-provider relay	SKIP [†]	—

[†] Skipped: requires three distinct backends; only two were available.

chosen deliberately: GPT-4o represents frontier-class capabilities while Llama 3.2 1B represents the smallest practical model, maximizing the diversity of the test.

Table 7 summarizes the cross-provider memory transfer tests. Each test writes memories with Provider A, then verifies recall with Provider B in a fresh session sharing only the .amem file.

Table 8 reports provider switch tests, which simulate realistic deployment scenarios such as switching cloud providers or transitioning between local and hosted models.

Table 9 reports brain file integrity tests verifying the binary format’s provider-neutrality at the byte level.

Summary. Of 22 tests, 21 passed and 1 was skipped due to an unavailable third backend—a 100% pass rate on all executable tests. The total validation time was 796 seconds (13.3 minutes).

6.3 Model Capability Analysis

The cross-provider validation revealed a significant capability gap between frontier and small models that merits discussion.

Structured output reliability. GPT-4o produces well-formed JSON from `chat_json()` prompts with near-perfect consistency (>99% of calls). Llama 3.2 1B, by contrast, frequently produces malformed output: responses wrapped in markdown code fences (````json...````), truncated JSON with missing closing braces, Python tuple syntax instead of JSON strings, invented schema keys not present in the prompt, and content mixed with explanatory prose.

Memory comprehension. Despite poor structured output, Llama 3.2 1B demonstrated strong *reading* capability: when provided with memory context retrieved from the .amem file, it correctly recalled facts, followed

Table 9: Brain file integrity tests confirming the .amem format contains no provider-specific data.

Test	Result	Time
No provider fingerprints (GPT-4o)	PASS	104.3 s
No provider fingerprints (Llama 1B)	PASS	86.1 s
File size sanity (GPT-4o)	PASS	36.3 s
File size sanity (Llama 1B)	PASS	96.8 s
Binary format consistency	PASS	4.9 s
Round-trip fidelity	PASS	4.8 s
Multi-provider brain health	PASS	27.4 s

correction chains, and maintained conversational continuity with the same accuracy as GPT-4o. This asymmetry—poor at structured *writing* but competent at unstructured *reading*—informed our layered robustness design.

Effectiveness of fallback layers. Before implementing the three-layer defense (Section 5.3), 4 of 22 cross-provider tests failed—all in scenarios where Llama 3.2 1B was the *writing* provider, producing JSON that the parser could not recover. After hardening, all tests passed. The regex fallback was activated on a significant fraction of Llama extraction calls, successfully recovering facts that would otherwise have been lost.

This finding has practical implications: AgenticMemory’s provider-agnostic design enables deployment on edge devices with small local models, provided the extraction pipeline accommodates their output variability. The three-layer approach achieves this without requiring model-specific code paths.

7 Discussion

Enabled capabilities. AgenticMemory makes several operations practical that are impossible with flat memory systems. *Cross-agent knowledge transfer* is achieved by copying a single file; the receiving agent inherits the full reasoning history, not just extracted facts. *Reasoning chain reconstruction* via CAUSED_BY traversal enables an agent to explain its decisions with full provenance. *Self-correction with history* via SUPERSEDES chains preserves the complete correction record, enabling meta-cognitive analysis of an agent’s belief evolution. *Portable identity* means an agent’s memory is not locked to a specific LLM provider; switching from one model to another requires no migration.

Empirical confirmation of portability. The cross-provider validation (Section 6.2) provides strong empirical evidence for the portability claim. Memories written by a 175B+ parameter cloud model were read correctly by a 1.24B parameter local model, and vice versa, across all tested scenarios. The binary format scan confirmed zero provider-specific fingerprints—no API tokens, model identifiers, or vendor-specific metadata leaked into the .amem file. This result is significant for deployment flexibility: organizations can switch LLM providers, or run different models for different use cases, without any memory mi-

gration cost.

The structured output gap. Our model analysis (Section 6.3) reveals that the primary barrier to cross-provider portability is not the memory format but the LLM’s ability to produce structured output for memory *formation*. Memory *retrieval* works consistently across model sizes because it requires only natural language comprehension. Memory *formation* requires JSON output, where small models struggle. The three-layer robustness pipeline (Section 5.3) bridges this gap, but a more principled solution might involve constrained decoding [21] or fine-tuning small models specifically for memory extraction.

Limitations. Feature vector generation still requires an external embedding model; AgenticMemory stores vectors but does not generate them. The current implementation does not support concurrent writes; the format assumes a single writer with multiple readers via memory-mapped access. The pattern query involves a linear scan and would benefit from bitmap or inverted indexes for large graphs.

8 Conclusion

We have presented AgenticMemory, a binary graph format purpose-built for persistent AI agent memory. By modeling memory as a graph of typed cognitive events connected by semantic edges, AgenticMemory enables capabilities that flat memory systems cannot provide: reasoning chain reconstruction, causal impact analysis, self-correction with full provenance, and cross-agent knowledge transfer.

The implementation delivers strong performance in a compact footprint: 4,795 lines of Rust for the core engine, 9,911 lines of Python for the agent framework, a 1.4MB compiled binary, zero external dependencies, and sub-10ms query latency on 100K-node graphs. Memory-mapped node access at 370 ns demonstrates that the fixed-size binary format achieves its $O(1)$ random access design goal. The entire memory of an agent—potentially years of accumulated knowledge—resides in a single portable file.

Empirical validation across 97 unit tests, 6 single-provider validation protocols, and 22 cross-provider tests (spanning GPT-4o and Llama 3.2 1B) confirms that the format is genuinely provider-agnostic: memories written by any supported backend are readable by any other, with zero data loss. The robustness engineering required to achieve this—JSON sanitization, retry pipelines, and regex fallback extraction—reveals practical challenges in deploying structured memory with heterogeneous LLM backends, challenges that are likely general to any system requiring reliable structured output from diverse models.

AgenticMemory represents a shift in how we think about agent memory: not as a search index over past text, but as a navigable graph of structured cognitive events.

We believe this graph-first approach is essential for building agents that can reason about their own history, explain their decisions, and learn from their mistakes.

References

- [1] Yann Collet. LZ4: Extremely fast compression. 2013. BSD licensed, version 1.9+.
- [2] Hermann Ebbinghaus. Über das gedächtnis. 1885. English translation: Memory: A Contribution to Experimental Psychology, 1913.
- [3] Google LLC. Protocol buffers: Google’s data interchange format. <https://protobuf.dev/>, 2008. Accessed 2025.
- [4] Google LLC. FlatBuffers: Memory efficient serialization library. <https://flatbuffers.dev/>, 2014. Accessed 2025.
- [5] Brook Heisler and Jorge Aparicio. Criterion.rs: Statistics-driven micro-benchmarking in Rust. <https://github.com/bheisler/criterion.rs>, 2023. Accessed 2025.
- [6] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.
- [7] LangChain Inc. LangMem: Memory management for LangChain agents. <https://github.com/langchain-ai/langmem>, 2024. Accessed 2025.
- [8] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474, 2020.
- [9] Mem0 AI. Mem0: Memory layer for AI agents. <https://github.com/mem0ai/mem0>, 2024. Accessed 2025.
- [10] Meta AI. Llama 3.2: Open-weight models for edge and mobile devices. <https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/>, 2024. Accessed 2025.
- [11] Ollama Inc. Ollama: Run large language models locally. <https://ollama.com/>, 2024. Accessed 2025.
- [12] OpenAI. ChatGPT memory: Personalization through persistent memory. <https://openai.com/index/memory-and-new-controls-for-chatgpt/>, 2024. Accessed 2025.
- [13] OpenAI. GPT-4o: Multimodal AI model. <https://openai.com/index/hello-gpt-4o/>, 2024. Accessed 2025.

- [14] OpenClaw Contributors. OpenClaw: Open-source memory for AI assistants. <https://github.com/openclaw>, 2024. Accessed 2025.
- [15] Charles Packer, Joseph E Gonzalez, Dawn Song, Raymond J Mooney, et al. MemGPT: Towards LLMs as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.
- [16] Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–22, 2023.
- [17] Pinecone Systems Inc. Pinecone: Vector database for machine learning. <https://www.pinecone.io/>, 2024. Accessed 2025.
- [18] Kenton Varda. Cap’n Proto: Insanely fast data interchange format. <https://capnproto.org/>, 2013. Accessed 2025.
- [19] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [20] Weaviate BV. Weaviate: Open-source vector database. <https://weaviate.io/>, 2024. Accessed 2025.
- [21] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, et al. Judging LLM-as-a-judge with MT-Bench and chatbot arena. In *Advances in Neural Information Processing Systems*, volume 36, 2023.