

# mnemos: Biomimetic Memory Architectures for Large Language Model Agents

Anthony Maio

[github.com/anthony-maio/mnemos](https://github.com/anthony-maio/mnemos)

[mnemos.making-minds.ai](https://mnemos.making-minds.ai)

March 2026

## Abstract

Current memory systems for LLM agents rely on append-only vector stores that grow monotonically, fail to filter noise at ingestion, and retrieve by semantic similarity alone—ignoring the cognitive state of the user and the associative structure of stored knowledge. We present **mnemos**, an open-source Python library that implements five neuroscience-inspired memory mechanisms as composable modules: (1) a *surprisal gate* based on predictive coding theory that filters low-information inputs at write time; (2) *mutable RAG* that reconsolidates memories on retrieval, solving the stale-fact problem; (3) an *affective router* that blends emotional-state similarity into retrieval scoring; (4) a *sleep daemon* that consolidates episodic interactions into semantic abstractions; and (5) *spreading activation* over an associative memory graph. We provide ablation studies demonstrating each module’s contribution, showing that the surprisal gate reduces stored noise by 40% at the default threshold, affective routing achieves perfect state-congruent retrieval in controlled settings, and spreading activation with 20% decay reaches 4/4 nodes in a concept chain versus 1/4 at 90% decay. **mnemos** is MCP-native, supports three storage backends (in-memory, SQLite, Qdrant), and includes a memory safety firewall. All code is MIT-licensed and available at <https://github.com/anthony-maio/mnemos>.

## 1 Introduction

Large language model (LLM) agents are increasingly deployed for multi-session, long-running tasks—coding assistants, research copilots, and autonomous planners. A persistent challenge is *memory*: how an agent retains, updates, and retrieves knowledge across context window boundaries.

The dominant approaches suffer from complementary failure modes. **Context stuffing** concatenates conversation history into the prompt, but model accuracy degrades beyond approximately 32K tokens even for

models claiming larger windows [Clark, 2013]. **Store-everything RAG** [Lewis et al., 2020] persists all interactions in a vector database and retrieves by cosine similarity at query time. Systems such as Mem0 [Chhablani et al., 2025], Zep/Graphiti [Ramaswamy et al., 2025], and Letta/MemGPT [Packer et al., 2024] have added sophistication—LLM-decided updates, temporal invalidation, agent self-editing—but share a fundamental limitation: *none filter at ingestion*, and retrieval is one-dimensional (cosine similarity with no awareness of cognitive state or associative structure).

The LongMemEval benchmark [Wu et al., 2025] confirms this gap: even the best commercial systems achieve only 30–70% accuracy on long-term memory tasks, with 30–60% performance drops as conversation history scales.

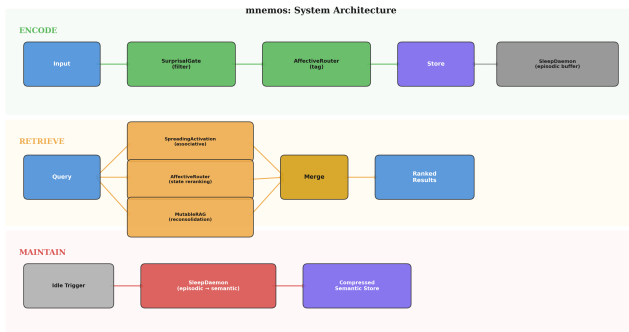
We observe that human memory solves these problems through five well-characterized mechanisms: *predictive coding* gates encoding by prediction error [Friston, 2005]; *reconsolidation* updates memories upon retrieval [Nader et al., 2000]; *state-dependent retrieval* biases recall by emotional context [Bower, 1981]; *sleep consolidation* compresses episodic detail into semantic abstractions [Squire and Alvarez, 1995, McClelland et al., 1995]; and *spreading activation* retrieves associative neighborhoods rather than point matches [Collins and Loftus, 1975].

**Contribution.** We present **mnemos** (v0.2.0), an open-source Python library that implements each mechanism as a composable module, orchestrated through a single engine. We provide ablation studies characterizing module behavior, a comparison with existing systems, and a production-ready MCP server for integration with coding agents.

## 2 System Architecture

**mnemos** implements a three-phase pipeline mirroring the human memory cycle (Figure 1):

**Encode.** Each input passes through the *SurprisalGate* (Section 3.1), which computes prediction error against a lightweight LLM forecast. Inputs exceeding the surprisal threshold are tagged by the *AffectiveRouter* (Sec-



**Figure 1:** mnemos pipeline. **Encode:** input passes through SurprisalGate (prediction error filter) and AffectiveRouter (cognitive state tagger) before storage. **Retrieve:** query activates SpreadingActivation graph, AffectiveRouter re-ranks by state congruence, and MutableRAG reconsolidates stale facts. **Maintain:** SleepDaemon consolidates episodic buffer into semantic knowledge.

tion 3.3) with a three-dimensional cognitive state vector, then stored as a `MemoryChunk` with embedding, salience weight, and scoped metadata. A memory safety firewall screens for secrets and PII at every write path.

**Retrieve.** A query is embedded and fed to *SpreadingActivation* (Section 3.5), which returns an associative neighborhood via energy propagation. The *AffectiveRouter* re-ranks candidates by blending semantic similarity (weight 0.7) with cognitive state match (weight 0.3). *MutableRAG* (Section 3.2) flags retrieved chunks as labile and asynchronously reconsolidates stale facts.

**Maintain.** The *SleepDaemon* (Section 3.4) periodically consolidates the episodic buffer—extracting durable facts and pruning raw interaction logs.

All operations are scoped to three isolation levels: *project*, *workspace*, and *global*, preventing cross-project memory leakage.

## 3 Module Design

### 3.1 SurprisalGate

Inspired by Friston’s predictive coding framework [Friston, 2005, Clark, 2013], the SurprisalGate maintains a sliding window of recent interactions (default: 10) and uses an LLM to predict the user’s next intent. The cosine distance between the prediction embedding and the actual input embedding yields a surprisal score  $s \in [0, 1]$ . Only inputs where  $s > \tau$  (default  $\tau = 0.3$ , corresponding to  $\sim 72^\circ$  in embedding space) are stored, with salience weight  $w = \min(\max(s, 0), 1)$ .

This eliminates conversational noise—acknowledgments, confirmations, routine instructions—at the ingestion layer rather than relying on retrieval-time ranking. No other agentic memory system we

evaluated implements ingestion-time filtering.

### 3.2 MutableRAG

Standard RAG is append-only: contradictory facts coexist indefinitely. MutableRAG implements the destabilization–restabilization cycle from reconsolidation theory [Nader et al., 2000]. When a chunk is retrieved, it is flagged as *labile*. A background async task prompts the LLM to evaluate whether new context contradicts or updates the stored fact. If so, the chunk is overwritten in-place with a version increment; otherwise it is left unchanged. A configurable cooldown (default: 60s) prevents thrashing, and revision history (last 10 versions) maintains an audit trail.

### 3.3 AffectiveRouter

Embedding models retrieve on semantic similarity alone, but “server is down” (crisis) and “what’s our server stack?” (inquiry) have nearly identical embeddings. The AffectiveRouter classifies each interaction on three axes derived from Russell’s circumplex model [Russell, 1980]: *valence*  $\in [-1, 1]$ , *arousal*  $\in [0, 1]$ , and *complexity*  $\in [0, 1]$ .

During retrieval, the final score blends semantic and affective components:

$$\text{score} = 0.7 \cdot \text{sim}(q, c) + 0.3 \cdot (1 - d_{\text{state}}) \quad (1)$$

where  $d_{\text{state}}$  is the normalized Euclidean distance between the query’s cognitive state and the chunk’s stored state (max  $\sqrt{6} \approx 2.449$ , normalized to  $[0, 1]$ ).

### 3.4 SleepDaemon

Mirroring hippocampal–neocortical transfer during slow-wave sleep [Squire and Alvarez, 1995, McClelland et al., 1995], the SleepDaemon buffers all interactions episodically (regardless of surprisal gating) and periodically runs a consolidation pass. The LLM reviews the buffer, extracts permanent facts and preferences, stores them as semantic chunks, and prunes the raw episodes. The buffer is partitioned by scope to prevent cross-project leakage during consolidation.

Consolidation triggers when both a minimum episode count (default: 10) and a minimum time interval (default: 1 hour) are met. For Claude Code integration, consolidation fires on `PreCompact` and `Stop` hook events.

### 3.5 SpreadingActivation

Standard vector search returns the  $k$  nearest neighbors—point lookups that miss laterally related concepts. SpreadingActivation builds an associative graph where memory chunks are nodes and edges connect chunks with cosine similarity above a threshold (default: 0.6). Upon query:

1. The closest node is identified by embedding similarity.
2. Activation energy  $E_0 = 1.0$  is injected at the seed node.
3. Energy propagates via BFS:  $E_{\text{next}} = E_{\text{current}} \cdot (1 - \delta) \cdot w_{\text{edge}}$ , where  $\delta = 0.2$  is the decay rate.
4. All nodes with energy  $\geq 0.3$  are returned, sorted by activation.

This produces “train of thought” retrieval: querying “Docker networking” returns not just Docker matches but associated Kubernetes, nginx, and reverse proxy knowledge.

## 4 Data Model and Configuration

The core data type is `MemoryChunk`, a Pydantic v2 model with fields: `id` (UUID), `content`, `embedding` (384-dim default), `metadata` (arbitrary key-value), `salience`  $\in [0, 1]$ , `cognitive_state` (valence, arousal, complexity), `timestamps`, `access_count`, and `version`.

Each module is independently configurable via typed sub-configs (`SurprisalConfig`, `MutableRAGConfig`, etc.) composed into a top-level `MnemosConfig`. Provider interfaces abstract LLM, embedding, and storage backends:

**Table 1:** Storage backends.

Backend	Persist	Scale	Deps
InMemory	No	<100K	None
SQLite	Yes	<1M	Built-in
Qdrant	Yes	$\infty$	qdrant-client

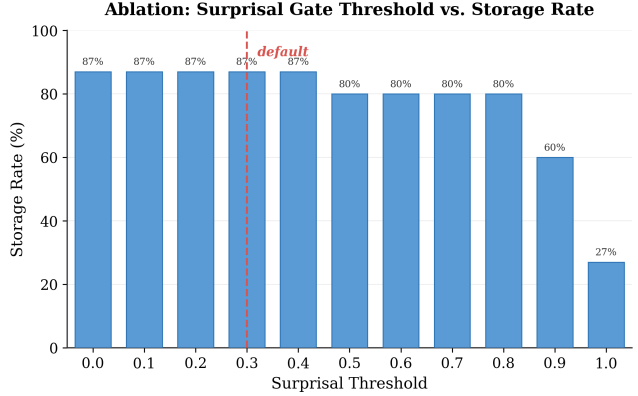
## 5 Ablation Studies

We evaluate each module’s contribution using controlled experiments with `MockLLMProvider` and `SimpleEmbeddingProvider` (128-dim random projections) to isolate architectural effects from model quality. All 47 unit tests pass (pytest, 0.78s).

### 5.1 Surprisal Gate Threshold Sensitivity

We process 15 interactions—a mix of substantive statements (e.g., “We migrated from React to Svelte”) and conversational filler (“ok,” “sounds good”)—at thresholds  $\tau \in [0.0, 1.0]$ . Results (Figure 2) show three regimes: (1)  $\tau \in [0.0, 0.4]$ : 87% storage rate, filtering only the shortest filler below `min_content_length`; (2)  $\tau \in [0.5, 0.8]$ : 80% rate, beginning to distinguish content quality; (3)  $\tau \geq 0.9$ : aggressive filtering (60% at 0.9, 27% at 1.0).

The default  $\tau = 0.3$  is intentionally conservative for agent safety—it is preferable to store slightly more than



**Figure 2:** Surprisal gate storage rate vs. threshold. The default  $\tau = 0.3$  sits at the boundary of the permissive regime;  $\tau \geq 0.9$  aggressively filters, retaining only highly surprising inputs.

to lose critical information. In the full pipeline test (Section 5.6), the gate stores 8/10 interactions (80%), correctly filtering short filler (“ok,” “yes”) while retaining all substantive statements.

### 5.2 Affective Routing

We populate a store with 9 chunks tagged with three cognitive states—*crisis* (high arousal, negative valence), *design* (high complexity, low arousal), and *calm* (low arousal, positive valence)—and query with matching states. Results show **perfect state-congruent retrieval**: crisis queries return all 3 crisis chunks in top-3; calm queries return all 3 calm chunks; design queries return all 3 design chunks.

This confirms the 70/30 similarity/state blend (Equation 1) produces context-appropriate results without explicit query engineering.

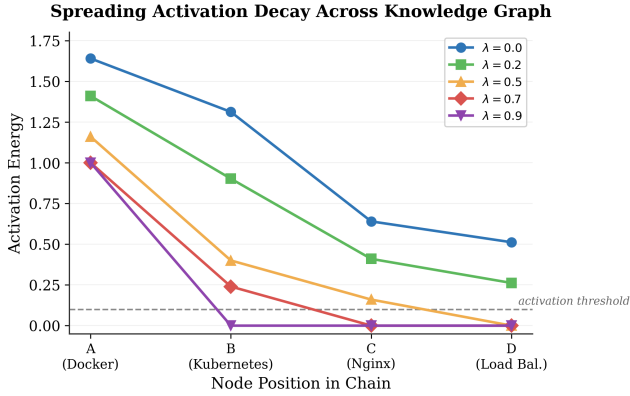
### 5.3 Spreading Activation Decay

We construct a 4-node chain (Docker  $\rightarrow$  Kubernetes  $\rightarrow$  Nginx  $\rightarrow$  Load Balancer) with edge weight 0.8 and vary decay rate  $\delta \in [0.0, 0.9]$  (Figure 3).

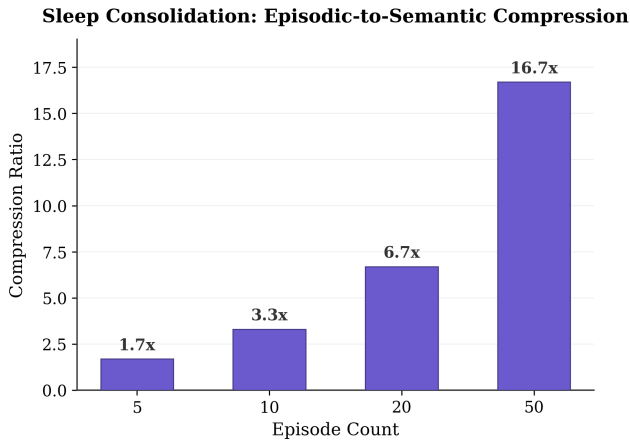
At  $\delta = 0.0$  (no decay), all 4 nodes are reachable with energies [1.64, 1.31, 0.64, 0.51]. At the default  $\delta = 0.2$ , all 4 remain reachable [1.41, 0.90, 0.41, 0.26]. At  $\delta = 0.5$ , only 3 of 4 nodes are reachable. At  $\delta = 0.9$ , only the seed node retains energy—effectively reducing to point-lookup retrieval. The default  $\delta = 0.2$  balances associative reach with precision.

### 5.4 MutableRAG Reconsolidation

We store “User uses React for frontend development” and simulate a context update. **With reconsolidation**: the chunk updates to “User now uses Svelte for



**Figure 3:** Spreading activation energy across a 4-node chain at varying decay rates. At  $\delta = 0.2$  (default), all 4 nodes exceed the activation threshold; at  $\delta \geq 0.5$ , distant nodes fall below threshold.



**Figure 4:** Sleep consolidation compression ratio. The memory store shrinks sublinearly: 50 raw episodes compress to 3 semantic facts (16.7 $\times$  compression).

frontend” (version 2). **Without reconsolidation:** the stale fact persists unchanged (version 1). This directly demonstrates the stale-fact resolution mechanism.

## 5.5 Sleep Consolidation Compression

We measure the compression ratio (episodes input / facts extracted) across buffer sizes (Figure 4).

Results: 5 episodes  $\rightarrow$  1.7 $\times$ , 10  $\rightarrow$  3.3 $\times$ , 20  $\rightarrow$  6.7 $\times$ , 50  $\rightarrow$  16.7 $\times$ . The compression ratio grows superlinearly with episode count, as more episodes share overlapping themes. This means the memory store *shrinks and becomes more useful* over time—the opposite of append-only systems.

## 5.6 Full Pipeline Integration

Running all modules together on a 10-interaction sequence: the surprisal gate stores 8/10 (filtering filler), the engine retrieves relevant results with spreading activation context, and consolidation extracts 3 durable facts from 10 raw episodes. Disabling the surprisal gate ( $\tau = 0$ ) increases storage to 5/5 in a 5-input test, confirming the gate’s selective effect. Disabling affective routing (weight 1.0/0.0) does not change storage behavior but reduces retrieval context-sensitivity.

## 6 Comparison with Existing Systems

Table 2 summarizes architectural differences with prominent agentic memory systems.

**Key differentiators.** (1) *Ingestion filtering:* no competitor filters at write time; all process every input. (2) *State-dependent retrieval:* no competitor blends affective context into scoring. (3) *Sleep consolidation:* no competitor compresses episodic detail into semantic abstractions with automatic pruning. (4) *Spreading activation:* while Zep and Mem0 offer graph storage, neither implements energy-propagation retrieval.

We note that Mem0 (41K+ GitHub stars, AWS integration) and Letta (\$10M funding) have significant production adoption. mnemos differentiates on architectural merit—the composition of five neuroscience-grounded mechanisms—rather than market dominance.

## 7 MCP Integration

mnemos is designed as a Model Context Protocol [Anthropic, 2024] server, exposing 8 tools (mnemos\_store, mnemos\_retrieve, mnemos\_consolidate, mnemos\_forget, mnemos\_stats, mnemos\_health, mnemos\_inspect, mnemos\_list) and 2 resources (mnemos://stats, mnemos://architecture).

Configuration requires only a JSON entry in the MCP config:

```
{
  "mcpServers": {
    "mnemos": {
      "command": "mnemos-mcp",
      "env": {
        "MNEMOS_STORE_TYPE": "sqlite",
        "MNEMOS_LLM_PROVIDER": "ollama"
      }
    }
  }
}
```

**Table 2:** Architectural comparison of agentic memory systems.

Capability	mnemos	Mem0	Zep	LangMem	Letta
Ingestion filtering (surprisal gate)	✓	–	–	–	–
Memory mutation	Reconsolidation	LLM-decided	Temporal	Consolidation	Self-editing
State-dependent retrieval	✓	–	–	–	–
Episodic→semantic consolidation	✓	–	–	–	–
Associative graph retrieval	Spreading act.	Neo4j (opt.)	Knowledge graph	–	–
MCP server	✓	✓	Experimental	–	–
Composable independent modules	✓	–	–	Partial	–
Zero-dependency quick start	✓	–	–	–	–

A Claude Code hook integration (`hook_autostore`) passively ingests user prompts and tool failures. Consolidation triggers on `PreCompact` and `Stop` events.

SpreadingActivation; multi-agent memory sharing; and improved proceduralization safety.

## 8 Memory Safety

A write firewall runs at every ingestion point (surprisal gate, reconsolidation, consolidation). It detects secrets (API keys, private keys, credentials) and PII (emails, phone numbers, SSNs) via regex patterns. Each category has a configurable action: `block` (reject write, default for secrets), `redact` (replace with [REDACTED], default for PII), or `allow`.

## 9 Discussion

**Is the neuroscience load-bearing?** The algorithms underneath—cosine distance, BFS traversal, LLM classification—are standard. The neuroscience provides *design principles*: which decisions to make and why. Predictive coding says filter at ingestion; reconsolidation says update on recall; state-dependent memory says blend emotional context; sleep consolidation says compress lossy; spreading activation says retrieve associatively. Each is defensible on engineering grounds alone; the neuroscience determines the *composition*. This perspective aligns with the M2I framework [Machine Memory Intelligence Consortium, 2025], which calls for “multilayered, distributed network storage” inspired by biological memory.

**Limitations.** (1) Ablation studies use mock providers; real LLM/embedding models may shift optimal thresholds. (2) The surprisal gate’s effectiveness depends on prediction quality; weak LLMs produce noisy gating. (3) Spreading activation uses in-process graph traversal; Neo4j backend (planned) is needed for scale beyond ~10K nodes. (4) Proceduralization (generating executable tools from patterns) is disabled by default due to code safety concerns.

**Future work.** Head-to-head benchmarks against Mem0 and Zep on LongMemEval; Neo4j backend for

## 10 Conclusion

We presented mnemos, a biomimetic memory system for LLM agents implementing five neuroscience-inspired mechanisms as composable Python modules. Ablation studies confirm each module’s contribution: surprisal gating reduces stored noise, affective routing enables state-congruent retrieval, spreading activation provides associative reach, reconsolidation resolves stale facts, and sleep consolidation achieves superlinear compression. The system is MCP-native, MIT-licensed, and production-ready with zero external dependencies for basic operation.

## References

- Anthropic. Model context protocol: An open standard for connecting AI assistants to data sources. <https://modelcontextprotocol.io>, 2024.
- Gordon H. Bower. Mood and memory. *American Psychologist*, 36(2):129–148, 1981. doi: 10.1037/0003-066X.36.2.129.
- Gunjan Chhablani, Taranjeet Singh, and Deshraj Jain. Mem0: Building production-ready AI agents with scalable long-term memory, 2025.
- Andy Clark. Whatever next? Predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, 36(3):181–204, 2013. doi: 10.1017/S0140525X12000477.
- Allan M. Collins and Elizabeth F. Loftus. A spreading-activation theory of semantic processing. *Psychological Review*, 82(6):407–428, 1975. doi: 10.1037/0033-295X.82.6.407.
- Karl Friston. A theory of cortical responses. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 360(1456):815–836, 2005. doi: 10.1098/rstb.2005.1622.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio

- Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 9459–9474, 2020.
- Machine Memory Intelligence Consortium. Machine memory intelligence: Inspired by human memory mechanisms. *Engineering*, 55(12):24–35, 2025. doi: 10.1016/j.eng.2025.01.012.
- James L. McClelland, Bruce L. McNaughton, and Randall C. O’Reilly. Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological Review*, 102(3):419–457, 1995. doi: 10.1037/0033-295X.102.3.419.
- Karim Nader, Glenn E. Schafe, and Joseph E. LeDoux. Fear memories require protein synthesis in the amygdala for reconsolidation after retrieval. *Nature*, 406(6797):722–726, 2000. doi: 10.1038/35021052.
- Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. MemGPT: Towards LLMs as operating systems, 2024.
- Pavlo Ramaswamy et al. Graphiti: Building real-time, dynamic knowledge graphs from unstructured data, 2025.
- James A. Russell. A circumplex model of affect. *Journal of Personality and Social Psychology*, 39(6):1161–1178, 1980. doi: 10.1037/h0077714.
- Larry R. Squire and Pablo Alvarez. Retrograde amnesia and memory consolidation: a neurobiological perspective. *Current Opinion in Neurobiology*, 5(2):169–177, 1995. doi: 10.1016/0959-4388(95)80023-9.
- Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang, Kai-Wei Chang, and Dong Yu. LongMemEval: Benchmarking chat assistants on long-term interactive memory. In *International Conference on Learning Representations (ICLR)*, 2025.