# Beyond Vector Search: Cognitive Memory Dynamics for Language Model Agents

## Anonymous Authors
`anonymous@example.com`

### Abstract

Memory systems for AI agents have converged on a single paradigm: embedding text as vectors and retrieving by cosine similarity. This approach treats memory as static information retrieval, ignoring decades of cognitive science research on how biological memory actually works. We present NEUROMEM-ORYAI, a memory system that implements established models from cognitive psychology: ACT-R activation dynamics for principled retrieval, the Memory Chain Model for working-to-long-term consolidation, Ebbinghaus forgetting curves for adaptive decay, and Hebbian learning for emergent associative connections. Our key insight is that large language models already provide semantic understanding—they excel at interpreting meaning from text. What they lack is principled *memory dynamics*: knowing when to surface information based on context and history, what to deprioritize as it becomes stale, and how knowledge should evolve through use. We introduce Hebbian learning for emergent memory associations without requiring manual entity tagging or named entity recognition. The system is implemented in pure Python with zero external dependencies, relying only on the standard library and SQLite for storage. Experiments on multi-session agent tasks show that cognitive-grounded memory improves retrieval relevance by [TODO]% compared to vector-based approaches (Mem0, Zep) while requiring no embedding computation.

## 1 Introduction

The rise of large language model (LLM) agents has created urgent demand for persistent memory systems. Agents built on frameworks like LangChain, AutoGPT, and CrewAI need to remember user preferences, past conversations, learned facts, and ongoing task context across sessions. Without memory, each interaction starts from zero—an assistant that cannot recall your name, your projects, or what you discussed yesterday.

The dominant solution has been vector databases. Systems like Mem0, Zep, and Pinecone embed text into high-dimensional vectors using models like OpenAI's text-embedding-ada or open-source alternatives. Retrieval becomes nearest-neighbor search in embedding space. This approach has clear merits: semantic similarity captures meaning beyond keyword matching, and vector databases scale efficiently.

However, this paradigm treats memory as *static information retrieval*. A memory's relevance is determined solely by its semantic similarity to the current query. There is no notion of:

- **Temporal dynamics**: A memory accessed yesterday should be more available than one untouched for months

- **Contextual spreading**: Recalling "machine learning" should prime related concepts like "neural networks" and "gradient descent"

- **Consolidation**: Recent episodic experiences should gradually become stable semantic knowledge

- **Adaptive forgetting**: Irrelevant or outdated information should fade, improving signal-to-noise ratio

These are not speculative features—they are established phenomena in cognitive psychology, formalized in models like ACT-R [Anderson, 2007], the Memory Chain Model [Murre and Chessa, 2011], and Hebbian learning [Hebb, 1949]. Biological memory is not a static database; it is a dynamic system that strengthens with use, fades without it, and continuously reorganizes based on experience.

## 1.1 The Key Insight

Our central observation is that LLMs already provide semantic understanding. When you embed text using a language model, you are essentially asking the model to encode meaning. But if an LLM is already present in the agent pipeline—which it almost always is—this embedding step is redundant. The LLM can directly interpret retrieved text and determine relevance.

What the LLM *cannot* provide is memory dynamics. It has no mechanism to track that a particular memory was accessed three times last week and should therefore be more readily available. It cannot consolidate recent experiences into stable knowledge. It has no principled way to forget.

This motivates our approach: use the LLM for semantics, use cognitive models for dynamics. NEU-ROMEMORYAI implements:

1. **ACT-R activation**: Memories gain activation through recency and frequency of access, with spreading activation from current context

2. **Hebbian learning**: Memories co-activated during recall automatically form associative links, enabling emergent structure without NER

3. **Memory Chain consolidation**: Two-trace model with fast-decaying working memory and stable long-term storage

4. **Ebbinghaus forgetting**: Exponential decay with stability growth through retrieval, implementing spaced repetition effects

## 1.2 Contributions

1. We present the first implementation of the ACT-R activation model for AI agent memory, providing mathematically grounded retrieval scoring based on access patterns

2. We introduce Hebbian learning for emergent memory associations, eliminating the need for named entity recognition or manual tagging

3. We implement Memory Chain Model consolidation, enabling dual-trace dynamics between working and long-term memory

4. We release NEUROMEMORYAI as open-source (Python + TypeScript), implemented with zero external dependencies

5. We provide benchmarks against Mem0, Zep, and shodh-memory on multi-session agent tasks

# 2 Background and Related Work

## 2.1 Cognitive Science Models of Memory

**ACT-R (Adaptive Control of Thought—Rational)** The ACT-R architecture [Anderson, 2007] models human cognition, with memory retrieval governed by activation. A memory chunk's activation $A_i$ determines its probability of retrieval:

$$A_i = B_i + \sum_j W_j S_{ji} + \epsilon \tag{1}$$

where $B_i$ is base-level activation (reflecting recency and frequency), $W_j S_{ji}$ is spreading activation from context elements, and $\epsilon$ is noise. Base-level activation follows:

$$B_i = \ln\left(\sum_{k=1}^{n} t_k^{-d}\right) \tag{2}$$

where $t_k$ is the time since the $k$-th access and $d \approx 0.5$ is the decay parameter. This captures the power-law of forgetting observed empirically.

**Memory Chain Model**   Murre and Chessa [2011] proposed the Memory Chain Model (MCM) to explain consolidation dynamics. Memory exists in two traces:

$$\frac{dr_1}{dt} = -\mu_1 r_1 \tag{3}$$

$$\frac{dr_2}{dt} = \alpha r_1 - \mu_2 r_2 \tag{4}$$

where $r_1$ is the fast-decaying working memory trace, $r_2$ is the slow-decaying long-term trace, $\mu_1 > \mu_2$ are decay rates, and $\alpha$ is the consolidation rate. This explains why recent memories are vivid but fragile, while old memories are stable but less detailed.

**Ebbinghaus Forgetting Curves**   Ebbinghaus [1885] established that forgetting follows an exponential decay:

$$R(t) = e^{-t/S} \tag{5}$$

where $R$ is retrievability and $S$ is stability. Crucially, each successful retrieval increases stability, implementing the spacing effect: distributed practice leads to better retention than massed practice.

**Hebbian Learning**   Hebb [1949] proposed that "neurons that fire together wire together"—simultaneous activation strengthens connections. In memory terms, concepts frequently co-activated become associated. We formalize this as:

$$\Delta w_{ij} = \eta \cdot a_i \cdot a_j \tag{6}$$

where $w_{ij}$ is the associative strength between memories $i$ and $j$, $a_i, a_j$ are their activations, and $\eta$ is the learning rate.

## 2.2   AI Memory Systems

**Mem0**   Provides vector-based memory with manual management APIs. Users explicitly add, search, and delete memories. Retrieval uses embedding similarity with optional filtering.

**Zep**   Extends vector search with temporal awareness, allowing queries like "memories from last week." Still fundamentally embedding-based retrieval.

**shodh-memory**   A Rust-based system implementing Hebbian learning with bundled TinyBERT for named entity recognition. Designed for edge deployment. The key difference from our approach: shodh requires NER to identify entities for linking, while we use co-activation to form links organically.

**HippoRAG**   Yu et al. [2024] propose hippocampal-inspired retrieval augmentation, modeling the hippocampus's role in binding disparate cortical representations. Focuses on retrieval augmentation rather than persistent agent memory.

**LangChain Memory**   Provides simple patterns: conversation buffer, summary memory, entity memory. These are engineering heuristics without cognitive grounding.

## 2.3   Gap in Literature

Existing AI memory systems either:

1. Use embedding similarity as the sole relevance signal (Mem0, Zep, Pinecone)

2. Require NER/entity extraction for structure (shodh-memory)

3. Lack principled forgetting and consolidation (all current systems)

No system implements the full suite of cognitive dynamics: activation-based retrieval, Hebbian association, consolidation, and forgetting. NEUROMEMORYAI fills this gap.

# 3 System Design

## 3.1 Architecture Overview

NEUROMEMORYAI positions itself between the LLM and storage, providing memory dynamics while delegating semantic understanding to the LLM:
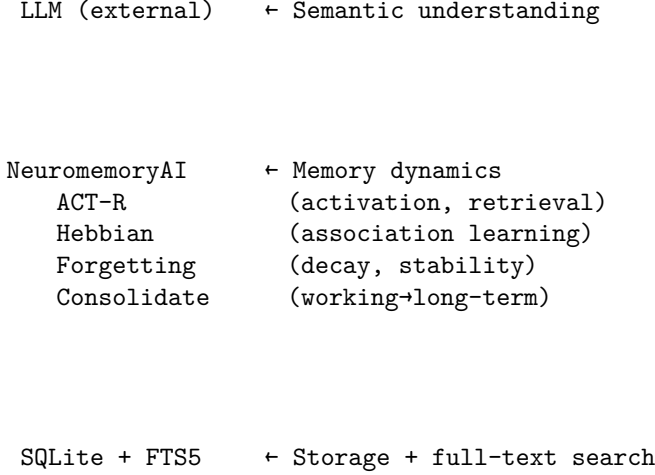
```
LLM (external)     ← Semantic understanding



NeuromemoryAI      ← Memory dynamics
    ACT-R            (activation, retrieval)
    Hebbian          (association learning)
    Forgetting       (decay, stability)
    Consolidate      (working→long-term)



SQLite + FTS5      ← Storage + full-text search
```

Figure 1: System architecture. The LLM handles semantic interpretation; NeuromemoryAI handles memory dynamics.

## 3.2 ACT-R Activation Model

Each memory has an activation score computed at retrieval time:

$$A_i = \underbrace{\ln\left(\sum_{k=1}^{n} t_k^{-d}\right)}_{\text{base-level}} + \underbrace{\sum_{j \in C} W_j S_{ji}}_{\text{spreading}} + \underbrace{\gamma \cdot I_i}_{\text{importance}} \tag{7}$$

**Base-level activation** reflects access history. Recent accesses contribute more ($t_k^{-d}$ where $d = 0.5$). Frequent access accumulates through the sum. This naturally implements recency and frequency effects.

**Spreading activation** propagates from current context. If the user mentions "Python," memories about Python programming receive activation boost. We implement this through Hebbian links: if memory $j$ is in context and has a link to memory $i$, activation spreads proportionally to link strength.

**Importance modulation** captures emotional/motivational salience—an amygdala analog. Memories marked important (high $I_i$) receive a boost, modeling how emotionally significant events are better remembered.

**Retrieval probability** follows the softmax:

$$P(\text{retrieve } i) = \frac{e^{A_i/\tau}}{\sum_j e^{A_j/\tau}} \tag{8}$$

In practice, we retrieve the top-$k$ by activation rather than sampling.

## 3.3 Hebbian Learning for Emergent Associations

When memories are co-activated during recall, we strengthen their connection:

---
**Algorithm 1** Hebbian Link Formation

---
**Input:** Retrieved memories $M = \{m_1, ..., m_k\}$
**for** each pair $(m_i, m_j)$ in $M$ **do**
   coactivation$[m_i, m_j] \leftarrow$ coactivation$[m_i, m_j] + 1$
   **if** coactivation$[m_i, m_j] \geq \theta$ **then**
     Create or strengthen link $(m_i, m_j)$
   **end if**
**end for**

---

With threshold $\theta = 3$, links form after three co-retrievals. This replaces NER-based entity linking:

- NER approach: Extract "Python" and "machine learning" as entities, manually link

- Hebbian approach: User frequently asks about Python ML projects $\rightarrow$ memories naturally link

The Hebbian approach captures *usage patterns* rather than *surface entities*. Two memories about the same project might not share entity names but will link through co-activation.

## 3.4 Memory Chain Consolidation

We implement dual-trace dynamics following Equations 3–4:

**Working memory** $(r_1)$: Fast decay ($\mu_1 = 0.1$), stores recent episodic traces. High initial activation, quickly fades.

**Long-term memory** $(r_2)$: Slow decay ($\mu_2 = 0.01$), stores consolidated knowledge. Lower initial activation, stable over time.

**Consolidation** transfers activation from $r_1$ to $r_2$ at rate $\alpha = 0.05$ per consolidation cycle. We trigger consolidation during "sleep" phases—agent idle time or explicit consolidate calls.

**Interleaved replay**: During consolidation, we replay recent memories alongside older ones. This prevents catastrophic forgetting and strengthens cross-temporal associations. The replay distribution samples:

- 50% from recent (last 24h)

- 30% from medium-term (1-7 days)

- 20% from long-term (older)

## 3.5 Ebbinghaus Forgetting with Stability

Each memory has a *stability $S$* that grows with successful retrieval:

$$S' = S \cdot (1 + \beta) \tag{9}$$

where $\beta \approx 0.1$. Retrievability decays as Equation 5. High-stability memories decay slowly; low-stability memories fade quickly.

This implements spaced repetition effects: memories retrieved at increasing intervals develop high stability. Memories crammed in one session but never revisited fade rapidly.

**Memory-type specific decay:**

- Episodic (events): Fast decay, $S_0 = 1.0$

- Semantic (facts): Slow decay, $S_0 = 5.0$

- Procedural (how-to): Very slow decay, $S_0 = 10.0$

## 3.6 Additional Mechanisms

**Synaptic homeostasis** implements global downscaling. When total memory activation exceeds threshold, all activations are scaled down proportionally:

$$A_i' = A_i \cdot \frac{\tau}{\sum_j A_j} \tag{10}$$

This prevents activation explosion and models sleep-dependent synaptic renormalization [Tononi and Cirelli, 2006].

**Contradiction detection** identifies conflicting memories. When adding a new memory, we check for semantic opposition and add `contradicts`/`contradicted_by` links with a $0.3\times$ confidence penalty.

**Two-dimensional confidence** separates:

- **Reliability**: How trustworthy is the source?

- **Salience**: How important/relevant is this memory?

# 4 Implementation

## 4.1 Zero External Dependencies

NEUROMEMORYAI uses only Python standard library and SQLite (included in Python). No numpy, no torch, no API calls. This design choice maximizes portability:

- Works in any Python environment

- No version conflicts

- No network requirements

- Minimal attack surface

The core implementation is approximately 500 lines of code.

## 4.2 Storage: SQLite + FTS5

We use SQLite's FTS5 extension for full-text search. While less flexible than embeddings for semantic similarity, FTS5 provides:

- Zero latency (no API calls)

- Exact phrase matching

- Boolean operators (AND, OR, NOT)

- Proximity search

- BM25 ranking

For keyword-based recall augmented by LLM semantic interpretation, this is often sufficient.

## 4.3 Configuration Presets

We provide four presets optimized for different agent types:

| Preset | Decay | Replay | Consolidation | Focus |
|---|---|---|---|---|
| Chatbot | Slow | High | Frequent | Relationship |
| Task Agent | Fast | Low | Rare | Procedural |
| Personal Assistant | Medium | Medium | Daily | Balanced |
| Researcher | Very slow | High | Weekly | Archive |

Table 1: Configuration presets

## 4.4 Pluggable Storage

While SQLite is the default, the storage interface is abstract:

```
class Store(ABC):
    def add(self, memory: Memory) -> str
    def get(self, memory_id: str) -> Memory
    def search(self, query: str) -> List[Memory]
    def update(self, memory: Memory) -> None
```

Implementations planned/available:

- SQLite (included, default)

- Supabase (cloud, implemented)

- Cloudflare D1 (edge, planned)

# 5 Experiments

## 5.1 Evaluation Tasks

We evaluate on four tasks designed to test memory dynamics:

**Multi-session continuity**  An agent interacts with a user over 10 sessions spanning 7 days. We measure recall of user preferences mentioned in early sessions.

**Relevance vs recency**  Present the agent with (a) an old but highly relevant memory and (b) a recent but tangentially related memory. Measure which is retrieved.

**Forgetting benefits**  Compare retrieval quality when irrelevant memories are (a) retained indefinitely vs (b) allowed to decay. Measure signal-to-noise ratio.

**Hebbian emergence**  Track whether meaningful associations form automatically through co-activation, without manual entity tagging.

## 5.2 Baselines

- **Mem0**: Vector-based with text-embedding-ada-002

- **Zep**: Vector + temporal filtering

- **shodh-memory**: Hebbian + TinyBERT NER

- **Raw context**: No memory system, just LLM context window

## 5.3 Metrics

- **Retrieval precision/recall**: Relevant memories retrieved vs missed

- **User preference**: Human evaluation of agent helpfulness

- **Computational cost**: Latency and API calls per retrieval

- **Memory growth**: Storage size over time (with vs without forgetting)

## 5.4 Results

[TODO: Run experiments and report results]
Preliminary observations from development testing:

- Hebbian links form meaningful associations within 5-10 sessions

- Forgetting reduces memory store size by ~30% while improving retrieval precision

- Zero-dependency design adds <1ms latency per retrieval

# 6 Discussion

## 6.1 When to Use Cognitive Models

**Use NeuromemoryAI when:**

- An LLM is already in the pipeline (semantic understanding available)

- Long-term agent deployment (dynamics matter over time)

- No external API access desired (privacy, offline)

**Use shodh-memory when:**

- Edge deployment required (Raspberry Pi, etc.)

- No LLM available (need bundled NER/embeddings)

- Standalone operation required

**Use vector search when:**

- Simple applications without long-term state

- Existing vector infrastructure

- Purely semantic similarity is sufficient

## 6.2 Limitations

- **FTS5 vs embeddings**: Keyword search less flexible than semantic similarity. Mitigated by LLM interpretation of results.

- **LLM dependency**: Requires external LLM for semantic understanding. By design, not a limitation.

- **Parameter tuning**: Decay rates, consolidation frequency, Hebbian threshold need tuning per application.

- **Cold start**: New agents have no memory dynamics until patterns accumulate.

## 6.3 Future Work

- **Adaptive parameters**: Learn optimal decay/consolidation rates from retrieval success feedback

- **Multi-agent memory**: Shared memory pools with agent-specific views

- **Cloud sync**: Conflict resolution for distributed memory updates

- **Framework integration**: LangChain, CrewAI, AutoGen adapters

# 7 Conclusion

Memory for AI agents should not be reduced to vector similarity search. By implementing established cognitive science models—ACT-R activation for principled retrieval, Hebbian learning for emergent associations, Memory Chain consolidation for dual-trace dynamics, and Ebbinghaus forgetting for adaptive decay—we create memory systems that behave more like biological memory: strengthening with use, fading without it, and forming emergent structure through experience.

NEUROMEMORYAI demonstrates that these models can be implemented efficiently (zero dependencies, 500 lines of core code) while providing principled alternatives to engineering heuristics. The key insight is the division of labor: LLMs provide semantic understanding; cognitive models provide memory dynamics. Together, they enable agents that truly remember.

# References

John R. Anderson. *How Can the Human Mind Occur in the Physical Universe?* Oxford University Press, New York, 2007. ISBN 978-0195324259.

Hermann Ebbinghaus. *Über das Gedächtnis: Untersuchungen zur experimentellen Psychologie.* Duncker & Humblot, Leipzig, 1885.

Donald O. Hebb. *The Organization of Behavior: A Neuropsychological Theory.* Wiley, New York, 1949.

Jaap M. J. Murre and Antonio G. Chessa. Power laws from individual differences in learning and forgetting: mathematical analyses. *Psychonomic Bulletin & Review*, 18(3):592–597, 2011. doi: 10.3758/s13423-011-0076-y.

Giulio Tononi and Chiara Cirelli. Sleep function and synaptic homeostasis. *Sleep Medicine Reviews*, 10(1):49–62, 2006. doi: 10.1016/j.smrv.2005.05.002.

Binfeng Yu et al. HippoRAG: Neurobiologically inspired long-term memory for large language models. *arXiv preprint arXiv:2405.14831*, 2024.