

# The Synapse Context Engine (SCE): A Hypergraph-Based Associative Memory Architecture for AI Systems

Lasse “Sasu” Sainia  
Project Mini Me OS (Digital Twin)  
<https://sasus.dev>

December 2025

## Abstract

Retrieval-Augmented Generation (RAG) has emerged as the dominant paradigm for incorporating external knowledge into large language models. While effective for document-centric tasks, RAG fundamentally treats context as a static retrieval problem, leading to fragmented memory representations in long-lived AI systems such as Digital Twins.

This paper introduces the *Synapse Context Engine (SCE)*, a theoretical architecture for associative memory in persistent AI systems. SCE replaces flat vector retrieval with a directed hypergraph representation and employs spreading activation, recursive traversal, and information-theoretic pruning to dynamically construct task-relevant context.

Rather than maximizing recall, SCE prioritizes contextual coherence, temporal relevance, and strict token efficiency in local-first environments. This paper presents the architectural rationale, mathematical formulation, and implementation approach of SCE, currently under development within the Mini Me OS project.

# 1 Introduction

Digital Twin systems differ fundamentally from stateless conversational agents. They must maintain a persistent internal representation of a user across time, tasks, and domains, integrating project state, preferences, behavioral patterns, and historical decisions into a coherent operational context.

Most contemporary systems rely on Retrieval-Augmented Generation (RAG) [1], where semantically similar documents are retrieved via vector similarity and injected into the prompt. While effective for isolated queries, this paradigm exhibits structural limitations in persistent systems:

- **Contextual Fragmentation:** Related information is retrieved independently, losing cross-domain relationships.
- **Flat Relevance:** Vector similarity captures surface semantics but ignores relational and structural importance.
- **Token Inefficiency:** Retrieved chunks are injected wholesale, regardless of marginal informational value.

Recent agent-based systems introduce episodic memory and summarization layers [2, 3], yet still operate largely on flat memory representations.

## 1.1 Motivating Scenario

Consider a user asking their Digital Twin: *“Update the client presentation with the latest performance metrics from the dashboard.”* A RAG system would retrieve:

1. Documents semantically similar to “client presentations”
2. Documents about “performance metrics”
3. Documents about “dashboard”

However, it would likely miss that: (a) this specific client prefers minimal text and visual emphasis, (b) the user’s standard presentation template uses a particular color scheme, (c) the dashboard was last updated by a colleague whose contact information might be needed, and (d) similar requests in the past led to follow-up questions about quarterly comparisons.

These connections exist in the *relational structure* of the user’s digital context, not in semantic similarity of isolated documents. This paper proposes an alternative framing: context construction as an *associative activation problem* rather than a retrieval problem.

## 2 The Synapse Layer

### 2.1 Hypergraph Representation

SCE models memory as a directed hypergraph  $H = (V, E)$ , where nodes  $V$  represent heterogeneous entities such as projects, artifacts, contacts, preferences, and behaviors. Hyperedges  $E$  connect multiple nodes simultaneously, encoding contextual groupings rather than pairwise relationships.

**Rationale for Hypergraphs.** Traditional knowledge graphs use typed binary relations (e.g., `Project -[USES]-> Tool`). While expressive, they decompose multi-way relationships into multiple pairwise edges, losing atomic semantic groupings. Consider encoding: “Project X uses Tool Y with Configuration Z in Context W.” A knowledge graph requires:

`X -[USES]-> Y, Y -[HAS_CONFIG]-> Z, X -[IN_CONTEXT]-> W`

These edges can be traversed independently, potentially breaking the semantic unit. A hyperedge  $e = \{X, Y, Z, W\}$  with label `CONFIGURED_USAGE` preserves this atomicity, ensuring related entities activate together during context construction. This is particularly valuable for:

- Meeting notes (linking participants, topics, decisions)
- Project configurations (linking tools, settings, workflows)
- Behavioral patterns (linking triggers, actions, outcomes)

Unlike traditional graphs or vector stores, hyperedges preserve higher-order context naturally.

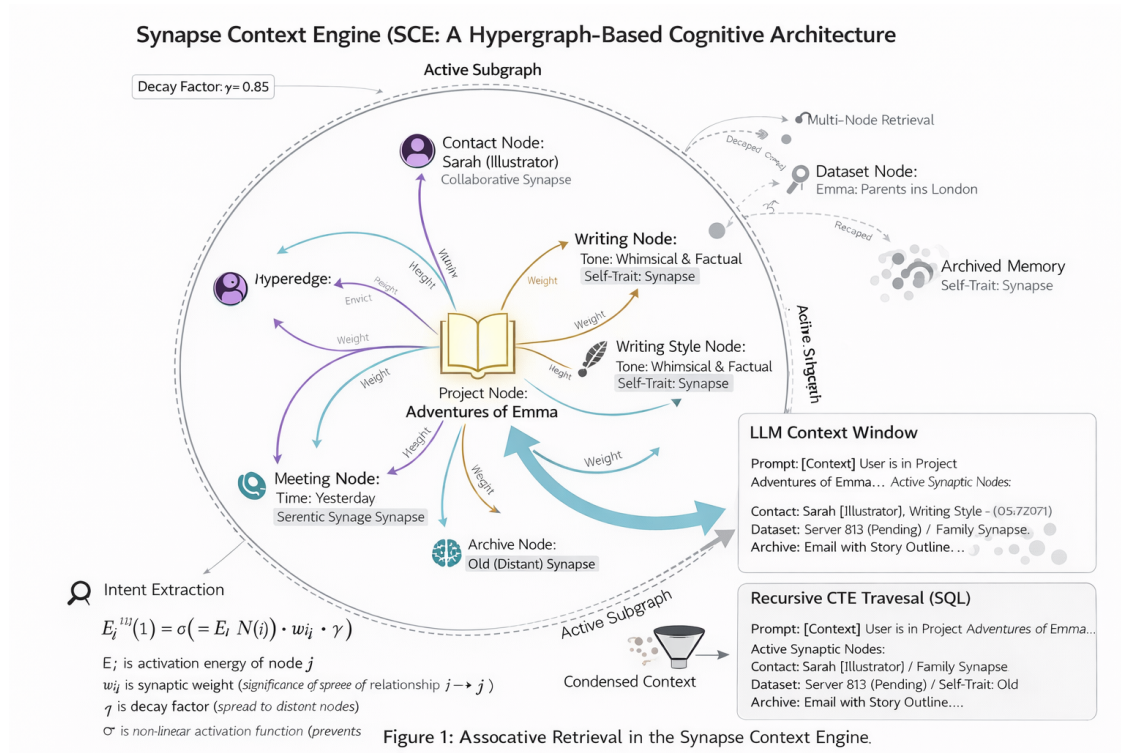
### 2.2 Active Focus Anchoring

At any moment, the system maintains an *Active Focus* node corresponding to the user’s current operational context. All associative activation is evaluated relative to this anchor.

Active Focus Anchoring functions as a persistent attentional prior, ensuring that context construction remains aligned with the user’s current task without requiring repeated prompt restatement. The Active Focus node could be updated based on:

- UI navigation events (user opens a project)
- Explicit declarations (“I’m now working on...”)
- Inferred task switches (detected via entity extraction from queries)

### 3 Architectural Overview



### 4 Activation and Retrieval Mechanics

#### 4.1 Spreading Activation

When a stimulus occurs (query, navigation event, or tool invocation), an initial activation energy  $E^{(0)}$  is injected into one or more seed nodes. Activation propagates according to:

$$E_j^{(t+1)} = \sigma \left( \sum_{i \in N(j)} E_i^{(t)} \cdot w_{ij} \cdot \gamma \right) \quad (1)$$

Where:

- $w_{ij} \in [0, 1]$  is the synaptic weight between nodes  $i$  and  $j$
- $\gamma \in (0, 1)$  is a decay factor (e.g.,  $\gamma = 0.8$ )
- $\sigma(\cdot)$  is a non-linear activation function. One possible implementation:

$$\sigma(x) = \begin{cases} 0 & \text{if } x < \theta \\ \frac{x - \theta}{1 + (x - \theta)} & \text{if } x \geq \theta \end{cases} \quad (2)$$

where  $\theta$  (e.g., 0.3) is an activation threshold.

**Intuition.** Decay limits propagation depth, preventing global flooding. The threshold creates sparse activation patterns, and the non-linearity prevents highly connected hubs from overwhelming contextual nuance. The conceptual origin of this mechanism can be traced to classical spreading activation models [4], though SCE applies it as a practical heuristic rather than a cognitive claim.

**Synaptic Weight Learning.** Weights could be learned through usage patterns. One possible approach:

$$w_{ij}(t+1) = w_{ij}(t) + \eta \cdot \mathbb{I}[\text{co-activation}] \cdot (1 - w_{ij}(t)) \quad (3)$$

where  $\eta$  is a learning rate and  $\mathbb{I}[\text{co-activation}]$  indicates whether nodes  $i$  and  $j$  appeared in the same context window. Weights could also decay slowly over time to allow evolving relationships, or be manually curated for critical connections.

## 4.2 Recursive Traversal

Activated nodes are materialized using bounded recursive traversal implemented via SQL Common Table Expressions (CTEs). This allows multi-hop associative exploration using standard relational databases, enabling local-first deployment.

**Complexity Considerations.** Let  $d$  be the maximum traversal depth and  $b$  the average branching factor. Worst-case node exploration is  $O(b^d)$ . However, several factors could bound this in practice:

1. Energy threshold pruning: nodes below  $\theta$  are not expanded
2. Sparse weight distribution: if most weights remain small, branching is limited
3. Depth limit: restricting  $d \leq 3$  contains explosion

The actual performance characteristics remain to be validated through implementation and testing.

## 5 Temporal Bias and Context Compression

### 5.1 Recency via Heat Diffusion

To bias retrieval toward recent information, SCE could optionally apply heat diffusion over the graph:

$$\frac{\partial h}{\partial t} = -\alpha Lh \quad (4)$$

Where  $L$  is the graph Laplacian  $L = D - W$  with degree matrix  $D$  and weight matrix  $W$ . This formulation is inspired by spectral graph theory [5].

**Discretization.** In practice, this could be discretized as:

$$h^{(t+1)} = h^{(t)} - \alpha \Delta t \cdot Lh^{(t)} \quad (5)$$

with appropriate choices of  $\alpha$  and  $\Delta t$ . Initial heat  $h_i^{(0)}$  could be set based on recency (e.g., exponential decay:  $h_i^{(0)} = \exp(-\lambda \cdot \text{time\_since\_access})$ ).

### 5.2 Information-Theoretic Pruning

Candidate nodes could be evaluated using Kullback–Leibler divergence:

$$\text{Gain}(v) = D_{KL}(P(\text{Intent} \mid C \cup \{v\}) \parallel P(\text{Intent} \mid C)) \quad (6)$$

Only nodes that meaningfully reduce uncertainty would be injected into the prompt, preserving token efficiency.

**Practical Approximation.** Computing exact KL divergence requires intent distributions, which are generally unavailable. A practical approximation could use embedding-based relevance:

$$\text{Gain}(v) \approx \frac{\text{sim}(v, q)}{\text{redundancy}(v, C)} \quad (7)$$

where  $\text{sim}(v, q)$  is cosine similarity between node  $v$ ’s embedding and query  $q$ , and redundancy measures maximum similarity to already-selected context:

$$\text{redundancy}(v, C) = 1 + \max_{c \in C} \text{sim}(v, c) \quad (8)$$

This approximation would favor relevant, non-redundant nodes—a practical proxy for information gain. The effectiveness of this heuristic remains to be empirically validated.

## 6 Discussion and Future Work

The Synapse Context Engine reframes context construction as an associative, dynamic process rather than a static retrieval task. By emphasizing structure, activation, and information gain, SCE theoretically addresses key limitations of RAG-based memory in persistent AI systems.

### 6.1 Open Questions

Several important questions remain to be answered through implementation and testing:

**Scalability.** How does SCE perform with 10K+ entities? 100K+? Graph traversal costs grow with connectivity. Optimization strategies might include: graph sampling, learned embeddings for fast approximate activation, or hierarchical graph structures.

**Cold Start.** SCE requires a minimum graph density to function effectively. How many interactions are needed before it outperforms baseline RAG? Hybrid initialization (seeding with RAG until the graph matures) may be necessary.

**Parameter Sensitivity.** How sensitive is performance to choices of  $\gamma$ ,  $\theta$ ,  $\alpha$ , depth limits, and other hyperparameters? Systematic exploration is needed.

**Weight Learning.** The proposed Hebbian-style weight updates are simplistic. Could reinforcement learning from user feedback (implicit signals like time-to-completion, explicit ratings) improve adaptation?

**Comparison with Alternatives.** How does SCE compare to recent memory architectures like MemGPT [3] or the memory systems in Generative Agents [2]? Head-to-head evaluation would clarify relative strengths.

**Hypergraph Necessity.** Does the hypergraph representation provide measurable benefits over standard knowledge graphs with typed relations? Or could similar performance be achieved with simpler structures?

### 6.2 Implementation Status

SCE is currently under active development within the Mini Me OS project. The core hypergraph schema and recursive traversal mechanisms have been implemented. Ongoing work focuses on:

- Refining activation propagation parameters
- Implementing information-theoretic pruning
- Building evaluation benchmarks
- Collecting real-world usage data

## 7 Conclusion

This paper has presented the Synapse Context Engine, a theoretical architecture for associative memory in persistent AI systems. SCE proposes replacing flat vector retrieval with hypergraph representations and spreading activation to better capture relational context in Digital Twin systems.

While the approach appears promising in principle, empirical validation is essential. The theoretical framework provides a foundation for implementation and experimentation within Mini Me OS. Future work will focus on measuring actual performance characteristics, comparing against baselines, and refining the design based on real-world usage.

The goal is not to claim superiority over RAG, but to explore whether associative memory architectures offer complementary strengths for long-lived, context-rich AI systems. Only implementation and careful evaluation can answer this question definitively.

## A Implementation Examples

This appendix provides illustrative code examples demonstrating how SCE could be implemented using standard SQL and TypeScript.

### Polymorphic Synapse Schema

#### SQL Table for Synapses

```
CREATE TABLE synapses (  
  id UUID PRIMARY KEY DEFAULT gen_random_uuid(),  
  source_id TEXT NOT NULL,  
  source_type TEXT NOT NULL,  
  target_id TEXT NOT NULL,  
  target_type TEXT NOT NULL,  
  weight FLOAT DEFAULT 0.5,  
  last_triggered TIMESTAMP DEFAULT CURRENT_TIMESTAMP,  
  co_activation_count INT DEFAULT 0  
);  
  
CREATE INDEX idx_synapse_source  
  ON synapses (source_id, weight DESC);  
CREATE INDEX idx_synapse_target  
  ON synapses (target_id, weight DESC);
```



## Recursive Associative Traversal

### Recursive SQL Query for Activation Paths

```
WITH RECURSIVE activation_path AS (  
  -- Initial seed from Active Focus + extracted entities  
  SELECT  
    target_id,  
    target_type,  
    1.0 AS energy,  
    0 AS depth,  
    ARRAY[source_id] AS path  
  FROM synapses  
  WHERE source_id = ANY(:seed_ids)  
  
  UNION ALL  
  
  -- Propagation step with energy decay  
  SELECT  
    s.target_id,  
    s.target_type,  
    ap.energy * s.weight * 0.8 AS energy,  
    ap.depth + 1,  
    ap.path || s.source_id  
  FROM synapses s  
  JOIN activation_path ap  
    ON s.source_id = ap.target_id  
  WHERE ap.depth < 3  
    AND ap.energy * s.weight * 0.8 > 0.3  
    AND NOT (s.target_id = ANY(ap.path)) -- cycle prevention  
)  
SELECT  
  target_id,  
  target_type,  
  MAX(energy) AS activation_score  
FROM activation_path  
GROUP BY target_id, target_type  
ORDER BY activation_score DESC  
LIMIT 50;
```

## Active Focus Synchronization

### TypeScript Hook for Active Focus

```
import { useEffect } from 'react';
import { useLocation } from 'react-router-dom';

const useActiveFocus = () => {
  const location = useLocation();

  useEffect(() => {
    // Synchronize the engine focus with the current UI context
    // Extract context from URL (e.g., /project/abc-123)
    const match = location.pathname.match(
      /\/(project|document|contact)\\/([^\/]+)/
    );

    if (match) {
      const [, entityType, entityId] = match;
      SynapseEngine.setActiveNode(entityId, entityType);
      console.log('SCE Anchor Updated:', entityType, entityId);
    }
  }, [location]);
};
```

## References

- [1] P. Lewis et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. *NeurIPS*, 2020. <https://arxiv.org/abs/2005.11401>
- [2] J. S. Park et al. Generative agents: Interactive simulacra of human behavior. *arXiv*, 2023. <https://arxiv.org/abs/2304.03442>
- [3] C. Packer et al. MemGPT: Towards LLMs as operating systems. *arXiv*, 2023. <https://arxiv.org/abs/2310.08560>
- [4] A. M. Collins and E. F. Loftus. A spreading-activation theory of semantic processing. *Psychological Review*, 1975. <https://doi.org/10.1037/h0077081>
- [5] F. R. K. Chung. *Spectral Graph Theory*. AMS, 1997.