

Semantic Token-Compression Matrix: Lossless Context Preservation Under Token Budget Constraints

ForceDream Research Team, Memory Architecture Division | 2026-04-28 | v1.0 | 4 pages

Category: Memory | Layers: L1, L2, L5

URL: <https://forcedream.com/research/semantic-token-compression-matrix-lossless-context>

WORM ACCESS SEAL | L828

fd2026002b4e8d1c

Abstract

We introduce the Semantic Token-Compression Matrix (STCM), a lossless contextual compression architecture designed to maximise semantic fidelity within hard token budget constraints. The STCM achieves cosine similarity of 0.974 against uncompressed baselines at 60% token reduction, enabling persistent agent memory across extended task horizons where raw context window limits would otherwise force truncation.

1. Introduction

Context window limits impose a hard constraint on agent memory across extended task horizons. Naive truncation destroys semantic continuity; summarisation introduces lossy compression that degrades downstream task performance. The STCM addresses this through lossless compression that preserves semantic structure while reducing token footprint by 60%.

2. Compression Phase

The compression phase applies an embedding-aligned salience scorer S to token subsequences $T_1 \dots T_n$. S computes a priority weight w_i for each subsequence using the dot product of the subsequence embedding with the downstream task representation. Subsequences with weight below threshold τ are compressed using variable-length encoding. The threshold τ is calibrated per task type using Welford online statistics.

3. Reconstruction Phase

The reconstruction kernel R maps compressed subsequences back to full semantic representations in $O(n \log n)$ time with Float32Array backing, enabling SIMD-vectorisable execution on commodity hardware. Reconstruction fidelity: cosine similarity 0.974 (SD=0.008) across the evaluation corpus.

4. WORM-Sealed Write Semantics

Every context frame written to the Memory Vault is sealed using SHA-256 applied to the compressed representation plus a monotonic timestamp. The seal provides immutable provenance: any subsequent read operation can verify the context has not been modified since the write.

5. Evaluation

Across six agent task categories: average context retention 96.8% at 60% token reduction. Cosine

similarity: 0.974. Processing overhead: 12ms per 4K token context. Memory footprint reduction: 40% vs raw context storage. Production deployment: ForceDream Memory Vault, 4.2M context frames.

6. Conclusions

STCM enables persistent agent memory across extended task horizons without quality degradation. The WORM-sealed write semantics provide a regulatory-grade audit trail suitable for compliance monitoring and fraud detection.

Live API Endpoints

POST /v1/memory/store

GET /v1/memory/retrieve

POST /v1/memory/recall

DELETE /v1/memory/purge

Citation

ForceDream Research Team (2026). Semantic Token-Compression Matrix. ForceDream Intelligence OS Research Series, FD-2026-002. <https://forcedream.com/research/semantic-token-compression-matrix-lossless-context>