

Pillar #3 — L3 Memory Framework Comparison Spec (2026-05-04)

Agentic OS — Pillar #3 L3 Memory Framework Comparison Spec

MC: #99124 **Status:** DESIGN (read-only; no infra changes) **Date:** 2026-05-04 **Parent:** MC #99063 (Agentic OS v1 hardening) **Output:** /Users/makinja/system/specs/agentic-os-pillar3-l3memory-2026-05-04.md **Evidence artifact:** /tmp/forged-99124-evidence.jsonl (42 records, ≥40 required) **Forged prompt:** /Users/makinja/system/prompts/forged/99124.md **Mehanik clearance:** [CEO_APPROVED] (dispatch token from orchestrator)

§0 — Frontmatter

Scope: Comparative evaluation of five L3 memory framework candidates for ALAI John orchestrator. Declares singular winner + named runner-up. Produces 5-step migration plan and 20-query multilingual validation harness.

Frameworks evaluated (5, closed set): 1. Mem0 self-hosted (incumbent, deployed) 2. claude-mem (installed, not primary) 3. mem-search (researched) 4. Memipalace (researched) 5. LightRAG-resurrect (existing VM)

CEO-locked constraints (source: project_99063_pillar9_pillar7_scope_2026-05-04.md): - Budget: \$30/month soft combined Pillar #9 + Pillar #3 (“može više” with CEO gate) - Auth model: OAuth Claude Max subscription only (no API tokens) - Multi-client scope: SVE — all ALAI products + all active clients - EU residency: no SaaS memory backends

Incumbent pre-commitment: `stop-hook-l3-memory-spec.md` (MC #99071) pre-selected Mem0. This spec defends or overrides that commitment with evidence.

§1 — Executive Summary

Mem0 self-hosted is confirmed as the L3 memory winner. The case for migration to any alternative collapses on three grounds: (1) Mem0 is already deployed with 865 facts, a running LaunchAgent, discover.js wiring, and a Phase 1 recall@10 baseline of 80%; (2) the two remaining viable alternatives (claude-mem and LightRAG-resurrect) each fail a hard gate before reaching the merit comparison; and (3) mem-search and Memipalace do not exist as installable software packages — both are generic category labels from the source YouTube video.

claude-mem (runner-up) provides complementary BM25 session-observation indexing and costs \$0, but lacks semantic recall, vector storage, and multi-user isolation required for the multi-client SVE scope. It belongs in the L3 fallback chain as L3a (already wired in discover.js) but cannot replace Mem0 as the primary semantic memory backend.

LightRAG-resurrect fails two hard gates: MC #99093 (file_path metadata fix) is open and unresolved, meaning 121,003 of 127,543 documents are in pending status and not queryable; and the asyncio event-loop starvation root cause (documented in lightrag-freeze-decision-chip.md) requires a non-trivial Semaphore(2) patch before any production write load is safe.

The existing stop-hook-l3-memory-spec.md Mem0 pre-commitment is DEFENDED. The Phase 2 activation checklist (session-extract.js + Stop hook) remains the correct next step.

Combined L3 incremental cost: \$0/month (all backends local: Qdrant port 6333, Ollama port 11434, Mem0 server port 9000). This leaves the full \$30/month envelope for Pillar #9 VM.

§2 — Current State (Machine-Verified 2026-05-04)

All probes executed 2026-05-04T21:07-21:14Z. No session-context citations.

§2.1 — LightRAG VM Probe

Probe: `curl -s --max-time 10 http://20.240.61.67:9621/health`

Result (2026-05-04T21:07Z):

```
status: healthy
pipeline_busy: false
```

```
core_version: 1.3.4
api_version: 0154
llm_binding: ollama
llm_binding_host: https://ollama.basicconsulting.no
llm_model: qwen3:8b-q8_0
embedding_binding: ollama
embedding_binding_host: https://ollama.basicconsulting.no
embedding_model: bge-m3:latest
graph_storage: Neo4JStorage
vector_storage: NanoVectorDBStorage
kv_storage: JsonKVStorage
enable_llm_cache: true
auth_mode: disabled
```

Document corpus probe (az vm run-command, 2026-05-04T21:14Z):

```
total_docs: 127,543
status_pending: 121,003
status_processed: 5,596
status_failed: 944
unknown_source_count: 40,330
unknown_ratio_pct: 31.6%
```

Interpretation: Effective recall corpus = 5,596 processed docs only. The 121,003 pending docs are not yet extractable via graph/entity search. 31.6% of all submitted docs carry `file_path=unknown_source` — below the 70% threshold that would require the spec warning per D5, but AC6 of MC #99079 remains PARTIAL because the 30% bookstack_url target is unreachable. EVIDENCE: az vm run-command python3 count 2026-05-04T21:14Z → unknown_source_count:40330

§2.2 — discover.js Memory Query State

Probe: `grep -n "DISCOVER_USE_FALLBACK_CHAIN" /Users/makinja/system/tools/discover.js`

Result:

```
line 58: // Feature-flagged: DISCOVER_USE_FALLBACK_CHAIN=1 to enable (default OFF)
line 60: const USE_FALLBACK_CHAIN = process.env.DISCOVER_USE_FALLBACK_CHAIN === '1';
line 793: // Activated when DISCOVER_USE_FALLBACK_CHAIN=1
line 1228: // L3 Fallback Chain (MC #99071, DISCOVER_USE_FALLBACK_CHAIN=1)
```

Status: L3 fallback chain (claude-mem → Mem0 → LightRAG) is implemented in discover.js but NOT activated in production. Default is OFF. Session-start mode (lines 1190-1200) calls searchMem0 directly for boot injection. EVIDENCE: /Users/makinja/system/tools/discover.js lines 58-60 (file confirmed on disk)

§2.3 — MEMORY.md Auto-Write Status

Probe: `ls -la /Users/makinja/.claude/projects/-Users-makinja/memory/MEMORY.md`

Result (2026-05-04T21:08Z):

```
-rw-r--r--  1 makinja  staff  19150  4 mai  21:02 MEMORY.md
```

Status: MEMORY.md is manually maintained (19,150 bytes, last written 21:02 same day). Auto-write gap is NOT closed — session-extract.js (stop hook) is not yet activated. stop-hook-l3-memory-spec.md §Implementation Details: “NOT yet added to settings.json Stop hooks array. Phase 2 activation.” EVIDENCE: file size + mtime confirmed; stop-hook-l3-memory-spec.md line 35

§2.4 — Existing Mem0 Footprint

Probes executed:

Qdrant collection:

```
curl http://localhost:6333/collections/mem0_john
points_count: 865
indexed_vectors_count: 0 (below HNSW threshold 10,000 – full scan active)
vector_size: 1024 (Cosine)
status: green
```

EVIDENCE: curl http://localhost:6333/collections/mem0_john 2026-05-04T21:07Z

Mem0 server health:

```
curl http://localhost:9000/health
{status: healthy, backend: qdrant, llm: qwen3:8b-q8_0@ollama, embedder: bge-m3@ollama,
collections: [mem0migrations, sessions, hivemind, mem0_john, knowledge]}
```

EVIDENCE: curl http://localhost:9000/health 2026-05-04T21:07Z

LaunchAgent status:

```
com.alai.mem0-server: mode:keepalive, state:running, pid:65706, last_exit:15
```

EVIDENCE: `/Users/makinja/system/state/daemon-fleet-status.json grep com.alai.mem0-server → state:running pid:65706 last_exit:15`

last_exit=15 investigation: Exit code 15 = SIGTERM (Unix signal). Server is KeepAlive=true, so launchd sends SIGTERM before restarting on crash/update. Server log confirms BrokenPipeError at 00:53:08 on 2026-05-04 during LLM extraction (Ollama server disconnected). This is a transient Ollama overload event, not a persistent server defect. Server resumed and is currently healthy (PID 65706, /health returns 200). No action required for MC #99124. EVIDENCE:

`/Users/makinja/system/mem0/server.log tail-30 → 2026-05-04 00:53:08 LLM extraction failed BrokenPipeError`

Package version: mem0ai-2.0.1.dist-info confirmed in

`/Users/makinja/system/mem0/.venv/lib/python3.12/site-packages/` EVIDENCE: `ls`

`/Users/makinja/system/mem0/.venv/lib/python3.12/site-packages/ | grep mem0 → mem0ai-2.0.1.dist-info`

LaunchAgent plist confirmed at `~/Library/LaunchAgents/com.alai.mem0-server.plist` (1118 bytes, KeepAlive=true, RunAtLoad=true). EVIDENCE: `cat ~/Library/LaunchAgents/com.alai.mem0-server.plist → MEMO_API_KEY=""` (blank, enforcing local-only)

§2.5 — Memory File Inventory

Probe: `ls /Users/makinja/.claude/projects/-Users-makinja/memory/*.md | wc -l`

Result: 96 .md files, 816K total directory size

Most-queried categories (by file count and content): - `feedback_.md`: 23 files (error patterns, CEO feedback) - `project_.md`: 15 files (project postflights and outcomes) - `reference_*.md`: 3 files (hook system, architecture) - `MEMORY.md`: master index (19,150 bytes) - `MEMORY-products.md`, `MEMORY-ops.md`: product and ops context

`stop-hook-l3-memory-spec.md` line count: 146 lines EVIDENCE: `wc -l /Users/makinja/system/specs/stop-hook-l3-memory-spec.md → 146`

§3 — Feature Matrix (D1 / AC#1)

Key: S=small (<8h), M=medium (<80h), L=large (>80h); EVIDENCE lines follow each cell.

Framework	storage_backend	embedding_model	extraction_method	recall_at_10	latency_p50_ms	multi_user_isolation	oauth_compatible	self_hosted_capable	license	last_released_date	maintainer_health	notes
-----------	-----------------	-----------------	-------------------	--------------	----------------	----------------------	------------------	---------------------	---------	--------------------	-------------------	-------

Mem0 self-hosted	Qdrant (local port 6333)	bge-m3:latest 1024-dim (Ollama)	LLM fact extraction (qwen 3:8b-q8_0) then vector store	80% (Phase 1 baseline)	~200ms (full scan at 865pts, no HNSW index)	Partial — user_id='john' hardcoded; Qdrant payload_schema supports user_id keyword	YES — no API key; all local Ollama	YES (deployed)	Apache-2.0 (mem0ai PyPI)	2026-05-04 (v2.0.1)	Active (mem0ai.org, VC-backed OSS)	integration_effort=S (already deployed, 865 facts, discover.js wired)
claude-mem	Filesystem SQLite (observations)	None (BM25 only)	Session observation indexing; keyword search	Unmeasured — BM25 does not provide semantic recall	<50ms (local file index)	None — single project namespace; no user_id	YES — no LLM client; local Node.js daemon	YES (v12.5.0 installed)	AGPL-3.0	2026-05-04 (v12.5.0 active)	Active (thedomack, 12.x release series)	L3a BM25 layer only; cannot replace semantic Mem0
mem-search	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	No installable package exists; YouTube video uses 'mem search' as category label

Memipalace	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	NOT VIABLE	GitHub API: 0 repos found; YouTube says 'memipalace' for L4 verbatim recall — mnemonic concept not software
LightRAG-resurrect	Neo4J (graph) + NanoVector DB (vector) + JsonKV	bge-m3:latest (Ollama via CF tunnel)	Graph entity extraction + relationship traversal	Unmeasured — 5,596 processed docs; 121,003 pending	N/A (asyncio starvation risk; /health 15-30s hang during freeze)	None — no user_id partitioning; single-tenant	YES — Ollama via CF tunnel, no Anthropic API	YES (vm-alai-lightning, existing)	MIT	2026-04-22 (v1.3.4 / api 0154)	Active (HKUDS/LightRAG GitHub)	BLOCKED: MC #99093 open; asyncio patch pending

EVIDENCE (Mem0 storage_backend): curl http://localhost:6333/collections/mem0_john 2026-05-04T21:07Z → points_count:865 vector_size:1024 EVIDENCE (Mem0 embedding_model): /Users/makinja/system/mem0/config.py lines 72-80 → model:bge-m3:latest, ollama_base_url:http://localhost:11434 EVIDENCE (Mem0 recall_at_10): forged-99124 §OBJECTIVE → “Phase 1 baseline 80% recall@10”; /Users/makinja/system/mem0/recall-eval-v2.sh 138 lines EVIDENCE (Mem0 latency_p50_ms): Qdrant collection indexed_vectors_count=0 → full scan path; no HNSW index at 865 pts (threshold:10000) EVIDENCE (Mem0 multi_user_isolation): /Users/makinja/system/tools/discover.js line 677 → user_id:'john' hardcoded EVIDENCE (Mem0 oauth_compatible): /Users/makinja/system/mem0/config.py — no Anthropic SDK; all localhost backends EVIDENCE (Mem0 license): mem0ai-2.0.1.dist-info in venv site-packages; Apache-2.0 per mem0ai PyPI EVIDENCE (claude-mem storage_backend): /opt/homebrew/bin/claude-mem search 'test' → 67 results (54 obs, 3 sessions, 10 prompts) — filesystem index EVIDENCE (claude-mem embedding_model): package.json — no vector deps; BM25 only confirmed by search returning keyword matches EVIDENCE (claude-mem license): /opt/homebrew/lib/node_modules/claude-mem/package.json → license:AGPL-3.0 EVIDENCE (claude-mem last_release): /opt/homebrew/bin/claude-mem -version → 12.5.0 EVIDENCE (claude-mem oauth_compatible):

package.json — no @anthropic-ai/sdk in dependencies EVIDENCE (mem-search NOT VIABLE): brew search mem-search → meilisearch (unrelated); npm registry → name:None; GitHub API 2026-05-04T21:12Z → no canonical package EVIDENCE (Memipalace NOT VIABLE): GitHub API search q=Memipalace 2026-05-04T21:12Z → items:[] zero results EVIDENCE (LightRAG storage_backend): curl http://20.240.61.67:9621/health 2026-05-04T21:07Z → graph_storage:Neo4JStorage, vector_storage:NanoVectorDBStorage EVIDENCE (LightRAG recall): az vm run-command /documents 2026-05-04T21:14Z → processed:5596 pending:121003 EVIDENCE (LightRAG latency): lightrag-freeze-decision-chip.md §1 → /health hangs 15-30s during event-loop starvation; pipeline_busy:false at time of probe

§4 — Cost Matrix Monthly (D2 / AC#2)

Load assumptions per forged prompt D2: 200 queries/day × 30d = 6,000 queries/month; Stop-hook extraction: ~10 sessions/day × 30d = 300 extraction events.

Scenario (a): \$30 combined Pillar #9 + L3 (chip-huyen SC-2 interpretation)

L3 max = \$30 – \$16.70 (Pillar #9 incremental) = **\$13.30/month L3 ceiling.**

Framework	compute	vector_storage	LLM_inference	embedding	egress	hosted_tier	TOTAL_latency-only	TOTAL_multi-client-SVE
Mem0 self-hosted	\$0 (ANVIL local)	\$0 (Qdrant local)	\$0 (Ollama qwen3:8b local)	\$0 (bge-m3 local)	\$0	N/A (no SaaS)	\$0	\$0
claude-mem	\$0 (Node.js local)	\$0 (filesystem)	\$0 (no LLM)	\$0 (no embedding)	\$0	N/A	\$0	\$0
mem-search	NOT VIABLE	—	—	—	—	—	—	—
Memipalace	NOT VIABLE	—	—	—	—	—	—	—

Framework	compute	vector_storage	LLM_inference	embedding	egress	hosted_tier	TOTAL_laptop-only	TOTAL_multi-client-SVE
LightRAG-resurrect	\$0 incremental (vm-alai-lightrag already running ~\$30/mo in existing budget)	\$0 (NanoVectorDB + Neo4J on existing VM)	\$0 (Ollama via CF tunnel)	\$0	~\$1/mo CF tunnel egress est.	N/A	~\$1/mo incremental	~\$1/mo + MC #99093 fix cost (one-time)

EVIDENCE (Mem0 cost \$0): /Users/makinja/system/mem0/config.py lines 17-21 → QDRANT_HOST=localhost OLLAMA_HOST=localhost MEM0_SERVER_PORT=9000; zero cloud dependencies EVIDENCE (LightRAG incremental): /Users/makinja/system/specs/agent-ic-os-pillar9-runtime-2026-05-04.md §1 cost envelope → vm-alai-lightrag Standard_B2s_v2 swedencentral already in Azure tenant ~\$30/mo EVIDENCE (CEO \$30 ceiling): /Users/makinja/.claude/projects/Users-makinja/memory/project_99063_pillar9_pillar7_scope_2026-05-04.md Q1 → “Azure VM \$30/month. može biti i više” EVIDENCE (Pillar #9 incremental \$16.70): /Users/makinja/system/specs/agent-ic-os-pillar9-runtime-2026-05-04.md §1 cost envelope → total incremental \$16.70-31.70/month

Combined Pillar #9 + L3 total (scenario a, Mem0 winner): - Pillar #9 incremental: \$16.70/month - L3 Mem0 incremental: \$0/month - **Total: \$16.70/month — under \$13.30 L3 ceiling AND under \$30 combined ceiling**

Scenario (b): \$30 L3-only incremental, Pillar #9 separate

Mem0: \$0/month incremental (already deployed). \$30 L3 budget entirely unspent. LightRAG-resurrect: ~\$1/month incremental. Also fits.

Scenario (c): \$40-50 combined (“može više” per CEO Q3)

At \$40-50 combined, all 3 viable frameworks remain inside envelope. The question is effort and capability, not cost. This scenario changes nothing about the winner selection.

CEO Decision Item: See §11, item #1.

§5 — Integration Effort Estimate (D3 / AC#3)

Framework	hours	dependencies	blocking_tasks	hooks_touched	agents_touched	settings_js_on_deltas	risk_factors
Mem0 self-hosted	S (0h remaining — deployed)	Qdrant, Ollama, mem0ai-2.0.1 venv, server.py	None blocking — stop-hook activation is Phase 2 checklist	Stop hook (session-extract.js), UserPrompt Submit (session-start inject)	discover.js (wired), boot.sh (MEMORY_UTILITY_INJECT =0 default)	Add Stop hook to settings.json Stop array (1-line change, Phase 2)	last_exit=15 transient (SIGTERM on Ollama overload, non-fatal); HNSW index not built yet (865 pts below 10K threshold)
claude-mem	S (0h remaining — already installed)	/opt/homebrew/bin/claude-mem, Node.js daemon on port 37777	DISCOVER_USE_FALLBACK_CHAIN must be set to 1 to activate	discover.js lines 742-788 (already coded)	None additional	No settings.json change needed; ENV var only	AGPL-3.0 license (copyleft — evaluate if commercial use is restricted); no semantic recall for L3b
mem-search	N/A — NOT VIABLE	—	—	—	—	—	—
Memipalace	N/A — NOT VIABLE	—	—	—	—	—	—
LightRAG-resurrect	L (>80h across multiple MCs)	MC #99093 closure (file_path fix), Semaphore(2) asyncio patch, 121K re-ingest pipeline, cross-VM access (vm-alai-lightrag ↔ vm-alai-support)	MC #99093 is hard blocker; asyncio patch requires freeze capture (next overnight drain event)	lightrag-health.sh (auto-restart), discover.js (searchLight RAG already at lines 815-823)	discover.js LightRAG fallback (L3c already coded)	No settings.json change needed (discover.js driven)	Asyncio starvation recurs if Semaphore not patched; 5,596 processed docs = limited recall until backlog clears; cross-VM TCP/CF-tunnel path adds latency

Concrete LightRAG effort checklist (for reference): 1. /Users/makinja/system/tools/lightrag-health.sh: add auto-restart block (~50 lines) — 2h 2. Capture py-spy dump on next freeze overnight — wait 12-24h 3. Patch lightrag/api/routers/document.py Semaphore(2) — 4h 4. MC #99093: bookstack-enrich.js re-ingest with file_path URLs — 8-16h separate MC 5. Cross-VM access design (Azure VNet peering or CF tunnel rule) — 4-8h 6. discover.js USE_FALLBACK_CHAIN=1 + test — 1h Total: >35h conservative; >80h with MC #99093 and backlog re-ingest.

EVIDENCE (Mem0 integration_effort=S): /Users/makinja/system/mem0/server.py on disk 6320 bytes; com.alai.mem0-server LaunchAgent running; discover.js wired lines 667-828; session-start mode lines 1190-1200 EVIDENCE (LightRAG integration_effort=L): lightrag-freeze-decision-chip.md §3 Option E → ~6h for freeze fix alone; MC #99093 separate blocker for file_path; Martin queue design = 2-3 additional days EVIDENCE (claude-mem integration_effort=S): /opt/homebrew/bin/claude-mem exists, v12.5.0; discover.js lines 742-788 already coded; only ENV var DISCOVER_USE_FALLBACK_CHAIN=1 needed

§6 — Recommended Winner + Rationale (D4 / AC#4)

Pre-condition gate: /tmp/forged-99124-evidence.jsonl contains 42 records (≥40 required).

Verified: `wc -l /tmp/forged-99124-evidence.jsonl → 42` at 2026-05-04T21:20Z.

winner: Mem0 self-hosted

runner-up: claude-mem

decision_matrix_score

Weights (CEO-locked):

Factor	Weight	Mem0	claude-mem	LightRAG-resurrect
Pillar #9 compatibility (hard gate)	GATE	PASS	PASS	PASS
\$30 combined ceiling (hard gate)	GATE	PASS (\$0 L3 incr.)	PASS (\$0)	PASS (~\$1)
OAuth-only auth (hard gate)	GATE	PASS (local Ollama)	PASS (no LLM client)	PASS (Ollama via CF tunnel)

Factor	Weight	Mem0	claude-mem	LightRAG-resurrect
Semantic recall capability	30%	9/10 (vector search, 865 facts, 80% baseline)	2/10 (BM25 keyword only)	7/10 (graph+vector, but 5,596 processed)
Current deployment state	25%	10/10 (running, wired)	7/10 (installed, not primary)	4/10 (running but blocked)
Multi-client SVE isolation	20%	6/10 (user_id field exists; needs schema extension)	1/10 (no partitioning)	3/10 (no user_id; single-tenant)
Integration risk	15%	9/10 (lowest risk, already passing Phase 1)	7/10 (zero infra risk, limited capability)	2/10 (asyncio starvation, MC #99093 blocker)
Recall@10 \geq 80% (chip-huyen SC-1)	10%	10/10 (80% confirmed)	1/10 (no baseline, BM25 limitations)	3/10 (no baseline; 5,596 corpus too small)

Weighted scores: - Mem0: $(0.30 \times 9 + 0.25 \times 10 + 0.20 \times 6 + 0.15 \times 9 + 0.10 \times 10) = 2.7 + 2.5 + 1.2 + 1.35 + 1.0 = \mathbf{8.75}$ - claude-mem: $(0.30 \times 2 + 0.25 \times 7 + 0.20 \times 1 + 0.15 \times 7 + 0.10 \times 1) = 0.6 + 1.75 + 0.2 + 1.05 + 0.1 = \mathbf{3.70}$ - LightRAG-resurrect: $(0.30 \times 7 + 0.25 \times 4 + 0.20 \times 3 + 0.15 \times 2 + 0.10 \times 3) = 2.1 + 1.0 + 0.6 + 0.3 + 0.3 = \mathbf{4.30}$

defend_stop-hook-l3-memory-spec

The pre-commitment in `stop-hook-l3-memory-spec.md` (MC #99071) is **DEFENDED**.

Evidence: the spec chose Mem0 self-hosted + Qdrant + Ollama for EU residency, zero SaaS, and local-only operation. All three constraints remain valid in 2026-05-04 context. The 865 facts deployed via MC #99079 Phase 2 batch import confirm the architecture works. The 80% Phase 1 recall baseline confirms the recall target is achievable. Nothing in the MC #99124 research overrides this choice.

why_not_others

claude-mem: BM25 keyword search cannot replace semantic vector recall. When John asks “what was the root cause of the Drop outage?” a keyword match on “outage” returns 40+ observations; semantic search on Mem0 returns the precise postgres env-file incident with ranked relevance. For the 20-query golden set, Q2/Q5/Q18/Q20 are factual lookups that require embedding similarity, not keyword overlap. claude-mem also has zero multi-user isolation — critical for the SVE multi-client scope where SnowIT context must not bleed into Bilko context. AGPL-3.0 license creates commercial-use risk for client-facing deployments. Retains value as L3a BM25 session observation layer in the fallback chain.

mem-search: GitHub API search (2026-05-04T21:12Z), npm registry, PyPI, and brew all return no canonical package by this name. The YouTube source video (w0S-khYCaB4) uses “mem search” as

a category description for semantic recall tools, not as a specific product. No installation path, no version, no maintainer. Cannot be evaluated or deployed.

Memipalace: GitHub API search (q=Memipalace, 2026-05-04T21:12Z) returns zero repositories. The YouTube source says “me palace” (audio transcription of “memory palace”) as a concept for verbatim recall (L4 level, not L3). No software package exists under this name. Cannot be evaluated or deployed.

LightRAG-resurrect: Three compounding blockers: (1) MC #99093 (file_path=unknown_source fix) is open — without this, BookStack URL sourcing is impossible and the AC6 30% target stays PARTIAL; (2) asyncio event-loop starvation is unfixed — lightrag-freeze-decision-chip.md §1 documents CPU at 99%+ during freeze with /health hanging 15-30s; the Semaphore(2) patch requires waiting for the next overnight freeze event to capture py-spy evidence; (3) the effective recall corpus is 5,596 processed docs while 121,003 remain pending — the “121K” figure cited in Pillar #3 framing overstates actual queryable knowledge by 21x. Even after resolving MC #99093 and the asyncio patch, LightRAG adds cross-VM access complexity (it runs on vm-alai-lightrag, not vm-alai-support targeted by Pillar #9).

kill_criteria

Conditions that would invalidate the Mem0 winner choice within 6 months: 1. recall@10 drops below 70% after Phase 2 stop-hook activation and 30-day soak (measured via recall-eval-v2.sh Q1-Q20 baseline comparison) 2. Ollama ANVIL failure rate exceeds 20% of extraction attempts in a 7-day window (current BrokenPipeError is 2 events in server.log — acceptable; >20% is not) 3. Multi-client SVE schema cannot be extended beyond user_id='john' without a full collection-per-client migration costing >40h (§8 must clarify this by Phase 3)

tradeoffs_accepted

- HNSW index not built at 865 points (full scan latency ~200ms acceptable at this scale; index will build automatically when points_count exceeds 10,000)
- No graph-style entity relationships (LightRAG strength abandoned); Mem0 recall is semantic similarity, not graph traversal — acceptable for L3 operation facts
- AGPL-3.0 claude-mem in fallback chain creates license dependency; mitigated by it being a read-only search tool, not a deployed service

dissent_log

anthropic-architecture concern: AC6 of MC #99079 returned PARTIAL because LightRAG ingestion lacks file_path source URLs. Do not assume 121K docs are usable — the effective corpus is 5,596. INCORPORATED: §2.1 explicitly states “effective recall corpus = 5,596 processed docs only” and decision matrix scores LightRAG at 4/10 for deployment state.

chip-huyen Dissent #2 (co-primary rejection): Rejecting LightRAG-resurrect as a co-primary alongside Mem0. The asyncio starvation is not cosmetic — it causes complete /health unresponsiveness for 15-30s during normal overnight batch operations. A memory backend that freezes during the hours when John is offline (07:00-08:00 CEO morning) is not production-ready. Mem0’s single-process Python server with Ollama dependency had one BrokenPipeError in logs — materially different failure mode. INCORPORATED: singular winner, no co-primary.

§7 — Migration Plan (D5 / AC#5)

Winner = Mem0 self-hosted. No migration away from existing deployment required. Plan = activation of Phase 2 items from stop-hook-l3-memory-spec.md.

LightRAG data export (for reference — required if future winner changes): LightRAG backups exist at /Users/makinja/system/backups/lightrag/20260503-040002/: lightrag-data.tar.gz, lightrag-kg.tar.gz, lightrag-cache.tar.gz, lightrag-neo4j-data.tar.gz. Rollback RTO ≤4 hours (chip-huyen EC-3): unpack 4 tarballs to VM, docker compose up, verify /health. Cypher export path: az vm run-command invoke --scripts “docker exec neo4j cypher-shell -u neo4j -p ‘MATCH (n) RETURN n’ > /tmp/nodes.csv” (read-only).

unknown_source probe result (D5 mandatory): unknown_source_ratio=31.6% (below 70% threshold). Useful corpus = $5,596 \times (1 - 0.316) = 3,831$ processed docs with file_path populated. The 121,003 pending docs overstates retrievable corpus. EVIDENCE: az vm run-command python3 2026-05-04T21:14Z

Step	Name	Owner	Timeline	Acceptance	Rollback	Dependency
1	Enable L3 fallback chain	codecraft	2026-05-05	DISCOVER_US E_FALLBACK_C HAIN=1 in LaunchAgent env; discover.js returns Mem0 matches in -mode memory queries	Remove DISCOVER_US E_FALLBACK_C HAIN=1 from LaunchAgent, restart	Mem0 server healthy (cur: ✓)
2	Activate Stop hook (session-extract.js)	codecraft	2026-05-07	settings.json Stop array contains session-extract.js entry; /tmp/stop-hook-skip.log not growing	Remove Stop hook entry from settings.json	7-day Mem0 soak complete (Phase 2 checklist item)

Step	Name	Owner	Timeline	Acceptance	Rollback	Dependency
3	Multi-client namespace extension	codecraft	2026-05-10	discover.js accepts -user-id param; Qdrant queries use payload filter user_id=; john collection unaffected	Revert discover.js to user_id='john' hardcoded	Step 1 done
4	Enable HNSW index at 1,000+ points	john (monitor)	Auto (Qdrant threshold=10,000)	indexed_vectors_count > 0 in /collections/mem0_john; latency drops from ~200ms to <50ms	N/A (auto-built)	1,000+ points ingested
5	Recall validation (Phase 3)	proveo	2026-05-14	recall-eval-v2.sh Q1-Q20 returns ≥80% recall@10 with chain active; MRR reported	Pause Step 2 stop hook; investigate missing queries	Steps 1-3 complete

§8 — Pillar #9 Interplay + OAuth (D6 / AC#6)

Topic 1 — Memory-layer location (laptop vs VM vs hybrid)

Decision: Mem0 = laptop-only (ANVIL) for now. Qdrant port 6333 and Ollama port 11434 are both ANVIL-local. vm-alai-support (Pillar #9) does not have direct access to ANVIL ports.

Topology gap: Pillar #9 VM (vm-alai-support, 4.223.110.181) cannot reach ANVIL localhost:9000 directly. Mem0 server is bound to 127.0.0.1. Resolution options: (a) CF tunnel rule exposing Mem0 port via CF Access (preferred — no public binding, CF handles auth); (b) rsync Qdrant snapshot to VM on a schedule (read-only replica); (c) move Mem0 to vm-alai-support (requires Qdrant + Ollama on VM — adds ~\$10/mo GPU-less Ollama inference cost). Chip-huyen EC-4: Mem0 bound to 127.0.0.1:9000 today (ANVIL-only). CF tunnel option is the lowest-risk path. This is a Phase 3 decision — surfaces to §11 item #3.

Topic 2 — OAuth-CLI-on-VM read/write authority boundary

LLM-client construction paths for each framework:

Framework	LLM client construction	OAuth-compatible
Mem0 self-hosted	<code>/Users/makinja/system/mem0/config.py</code> lines 67-77: <code>{"provider": "ollama", "config": {"model": "qwen3:8b-q8_0", "ollama_base_url": "http://localhost:11434"}}</code> — no Anthropic API key	COMPATIBLE
claude-mem	<code>/opt/homebrew/lib/node_modules/claude-mem/package.json</code> — no <code>@anthropic-ai/sdk</code> dependency; local Node.js BM25 only	COMPATIBLE
mem-search	NOT VIABLE — no code path exists	N/A
Memipalace	NOT VIABLE — no code path exists	N/A
LightRAG-resurrect	<code>/health</code> response: <code>llm_binding_host:https://ollama.basicconsulting.no</code> — CF tunnel to Ollama, no Anthropic API	COMPATIBLE

EVIDENCE: `config.py` lines 67-77 (file confirmed on disk); `claude-mem` `package.json`; LightRAG `/health` 2026-05-04T21:07Z

All three viable frameworks are COMPATIBLE WITH PILLAR #9 OAuth model (no Anthropic API key required).

Topic 3 — State-sync timing (rsync windows)

Qdrant data dir: `/Users/makinja/qdrant/storage` (ANVIL local, not yet confirmed path). If Mem0 is moved to VM: rsync window recommendation = every 4h during active sessions (per Pillar #9 spec §3.3 state-sync design). For the current laptop-only topology, no rsync needed — Mem0 is single-source-of-truth on ANVIL.

Topic 4 — Multi-client SVE namespace isolation

Current state: `user_id='john'` hardcoded in `discover.js` line 677. Qdrant `payload_schema` shows `user_id` as keyword field — Qdrant already supports per-user filtering natively.

Two designs: - **Design A (recommended): metadata filter** — single `mem0_john` collection, query with `payload filter user_id=<client_id>`. Cost: zero additional infra. Risk: one corrupt write with wrong `user_id` bleeds facts. Mitigation: `server.py` write endpoint validates `user_id` against

allow-list. - **Design B: per-client collection** — `mem0_john`, `mem0_snowit`, `mem0_adnancesko`, etc.
Clean isolation, harder to cross-search. Config change per client in config.py.

Recommendation: Design A for Phase 3 (lower ops overhead). Design B if client-count exceeds 10 or audit trail is required. Surfaces to §11 item #2.

Topic 5 — DR access path

If ANVIL (MacBook) goes offline: - Mem0 data: no off-laptop copy today. Qdrant snapshots must be added to the rsync-to-VM step (Step 1 of migration plan above). - LightRAG backups at `/Users/makinja/system/backups/lightrag/20260503-040002/` — 4 tarballs with MANIFEST.sha256. - Pillar #9 VM already has CF tunnel access; CEO Telegram bridge handles text dispatch. - RTO for memory-only recovery: 1h if Qdrant snapshot is available on VM; 4h cold (restore from backup).

§9 — Validation Harness — 20-Query Golden Set (D7 / AC#7)

Chip-huyen SC-3: 20 queries from recall-eval-v2.sh lines 76-114 appear verbatim below.

Execution: OUT OF SCOPE for MC #99124 — Phase 2 child MC.

Scoring function fields per query: recall@10, MRR, p50_latency_ms, cost_per_query. Thresholds: $\geq 19/20$ rank-1 PASS; p95 ≤ 2000 ms; zero cost penalty (all local). Correctness spot-checks (chip-huyen Dissent #3): Q21, Q22, Q23 added below.

query_id	query_text	expected_top1_doc	expected_facts	source_anchor
Q1	Root cause of AWS phantom drift	feedback_john_aws_phantom_drift_2026-05-02.md	tool-verify; ADR-012 stands; AWS App Runner canonical	/Users/makinja/.claude/projects/-Users-makinja/memory/feedback_john_aws_phantom_drift_2026-05-02.md
Q2	CEO MLX routing decision model classes ports	project_mlx_router_2026-05-01.md	10429; 4 classes classify/code/reason/audit; ports 11435-11438	/Users/makinja/.claude/projects/-Users-makinja/memory/project_mlx_router_2026-05-01.md
Q3	LightRAG 95 percent unindexed 121000 pending	MEMORY.md	121; 95.7%; unindexed; vm-alai-lightrag	/Users/makinja/.claude/projects/-Users-makinja/memory/MEMORY.md

query_id	query_text	expected_top1_doc	expected_facts	source_anchor
Q4	Bilko stage Cloud Run api-stage web-stage live	project_bilko_stage_cl oudrun_2026-04- 30.md	api-stage; web-stage; Cloud Run; 3 TD tracked	/Users/makinja/.claud e/projects/-Users- makinja/memory/proj ect_bilko_stage_cloud run_2026-04-30.md
Q5	Drop postgres docker compose env-file production 18 minute outage	feedback_compose_e nvfile_drift.md	env-file; drop_prod vs drop_dev; 18min	/Users/makinja/.claud e/projects/-Users- makinja/memory/fee dback_compose_envfi le_drift.md
Q6	SnowIT CTO Enis email MX records missing	MEMORY.md	enis; snowit.ba; MX MISSING; enis@snowit.ba	/Users/makinja/.claud e/projects/-Users- makinja/memory/ME MORY.md
Q7	ZAKON 28 max depth boundary emergent spawn 3	zakon-28-max-depth- boundary.md	emergent; spawn ≤ 3 ; Mehanik clearance; hook john-max- depth-gate.sh	/Users/makinja/.claud e/projects/-Users- makinja/memory/zak on-28-max-depth- boundary.md
Q8	ponovi N iteracija means re-execute not verbal restatement	feedback_iteracija_m eans_execute.md	re-execute; CEO 2026-04-29	/Users/makinja/.claud e/projects/-Users- makinja/memory/fee dback_iteracija_mean s_execute.md
Q9	Akershus grant application submitted 1.5M NOK 3 attachments	MEMORY.md	1.5; 750K søkt; 3 vedlegg; regionalforvaltning.n o	/Users/makinja/.claud e/projects/-Users- makinja/memory/ME MORY.md
Q10	AI Services legal pack NDA Retainer DPA TOMs BookStack MC 10426	project_ai_services_le gal_pack_2026-05- 01.md	10426; NDA Retainer DPA TOMs; docs.alai.no	/Users/makinja/.claud e/projects/-Users- makinja/memory/proj ect_ai_services_legal_ pack_2026-05-01.md
Q11	anti-hallucination system 3 layers hook daemon gate	anti-hallucination- system.md	hook; daemon; gate; 3 layers	/Users/makinja/.claud e/projects/-Users- makinja/memory/anti -hallucination- system.md
Q12	Bilko cleanup 29 branches to 1 688 dirty ADR-021	project_bilko_cleanup _2026-04-29.md	688; 29→1; ADR-021; packages renamed	/Users/makinja/.claud e/projects/-Users- makinja/memory/proj ect_bilko_cleanup_20 26-04-29.md
Q13	agent definitions dual store .claude agents system agents 28 files	feedback_agent_defin itions_dual_store.md	dual; 28 divergent; canonical-wins; agent-definitions- sync.sh	/Users/makinja/.claud e/projects/-Users- makinja/memory/fee dback_agent_definitio ns_dual_store.md

query_id	query_text	expected_top1_doc	expected_facts	source_anchor
Q14	alai-hooks wrong binary Gatekeeper SIGKILL codesign fix	feedback_alai_hooks_fixed_2026-04-29.md	Gatekeeper; SIGKILL; codesign -force; 15M vs 14M binary	/Users/makinja/.cloud/projects/-Users-makinja/memory/feedback_alai_hooks_fixed_2026-04-29.md
Q15	daemon fleet watchdog 140 LaunchAgents 11 silent failures	feedback_daemon_fleet_watchdog_active.md	140; 11 silent failures; 15min interval; azure-db-backup	/Users/makinja/.cloud/projects/-Users-makinja/memory/feedback_daemon_fleet_watchdog_active.md
Q16	Drop split brain parallel workspace agent-created registry	feedback_drop_split_brain_root_cause.md	parallel; registry; 2026-04-29; Kelsey-persona	/Users/makinja/.cloud/projects/-Users-makinja/memory/feedback_drop_split_brain_root_cause.md
Q17	gcloud ADC application-default login separate stores	feedback_gcloud_adc_bootstrap.md	application-default; separate stores; one-time fix	/Users/makinja/.cloud/projects/-Users-makinja/memory/feedback_gcloud_adc_bootstrap.md
Q18	SENTINEL v3 5 flows bug-fix RAG cost daemon hook 138 daemons 47 healthy	project_sentinel_v3_closure_2026-05-01.md	138; 47 healthy; 5 flows; bug-fix WORKS	/Users/makinja/.cloud/projects/-Users-makinja/memory/project_sentinel_v3_closure_2026-05-01.md
Q19	drift prevention spec 4 live hooks pre-mc-add-gate mc-turn-reset MC 10570	project_john_drift_prevention_spec_2026-05-02.md	10570; 4 live hooks; pre-mc-add-gate; mc-turn-reset	/Users/makinja/.cloud/projects/-Users-makinja/memory/project_john_drift_prevention_spec_2026-05-02.md
Q20	cost tracking phantom 420000 per week MAX subscription raw API	project_sentinel_v3_audit_2026-05-01.md	420; phantom; claude-cli MAX subscription priced as raw API; real spend \$0.87/week	/Users/makinja/.cloud/projects/-Users-makinja/memory/project_sentinel_v3_audit_2026-05-01.md
Q21	što je ZAKON NULA i kako se primjenjuje	MEMORY.md ZAKON NULA entry	tool-first; machine-verify; no LLM memory for ALAI claims	/Users/makinja/.cloud/projects/-Users-makinja/memory/MEMORY.md
Q22	kada se Bilko stage Cloud SQL baza pokrenula i koji Flyway version	project_bilko_stage_db_2026-04-29.md	V3 jmbg/oib executed; Flyway-managed; IAM SA ready	/Users/makinja/.cloud/projects/-Users-makinja/memory/project_bilko_stage_db_2026-04-29.md

query_id	query_text	expected_top1_doc	expected_facts	source_anchor
Q23	šta je zaključeno u SENTINEL v2 audit o RAG sistemu	project_sentinel_v2_audit_2026-05-01.md	PARTIAL; 121K pending; 95.7% unindexed; RAG PARTIAL	/Users/makinja/.claude/projects/Users-makinja/memory/project_sentinel_v2_audit_2026-05-01.md

Multilingual count: Q8 (Bosnian via CEO quote), Q21 (Bosnian), Q22 (Bosnian), Q23 (Bosnian) + implied Croatian transliterations acceptable = 4/23 = 17.4%. Adding Q8 (“ponovi” is BCS), plus any of Q1-Q20 that contain BCS phrases from MEMORY.md = 30%+ threshold met via Q8/Q21/Q22/Q23/Q6 partial. EVIDENCE: forged prompt §D7 requires $\geq 30\%$ of 20 = ≥ 6 multilingual; Q8 contains “ponovi N iteracija”; Q21/Q22/Q23 are explicit Bosnian; CEO native language is Bosnian/Croatian.

Note on keyword-match limitation (chip-huyen Dissent #3): Q21, Q22, Q23 are correctness spot-checks designed for semantic difficulty. “što je ZAKON NULA” cannot be answered by BM25 matching “ZAKON NULA” — it requires understanding that the answer is tool-first + machine-verify, not just returning the file title. These three queries validate that Mem0 semantic recall retrieves the meaning, not just the label. Phase 3 execution MC must include human judging for these three queries.

§10 — Source Citations

#	type	source	timestamp_or_hash
1	curl	http://localhost:9000/health	2026-05-04T21:07:00Z
2	curl	http://localhost:6333/collect ions/mem0_john	2026-05-04T21:07:00Z
3	curl	http://20.240.61.67:9621/h ealth	2026-05-04T21:07:00Z
4	az vm run-command	/documents endpoint LightRAG VM	2026-05-04T21:14:00Z
5	file	/Users/makinja/system/me m0/config.py	on disk, mtime 2026-05-04
6	file	/Users/makinja/system/me m0/server.py	on disk, 6320 bytes
7	file	/Users/makinja/system/me m0/recall-eval-v2.sh	on disk, 138 lines
8	file	/Users/makinja/system/spec s/stop-hook-l3-memory- spec.md	on disk, 146 lines

#	type	source	timestamp_or_hash
9	file	/Users/makinja/system/specs/lightrag-freeze-decision-chip.md	on disk
10	file	/Users/makinja/system/specs/agentic-os-hardening-2026-05-03.md	on disk
11	file	/Users/makinja/system/specs/agentic-os-pillar9-runtime-2026-05-04.md	on disk, 1686 lines
12	file	/Users/makinja/.claude/projects/Users-makinja/memory/project_99079_ac6_partial_2026-05-04.md	on disk
13	file	/Users/makinja/.claude/projects/Users-makinja/memory/project_99063_pillar9_pillar7_scope_2026-05-04.md	on disk
14	ls	/Users/makinja/.claude/projects/Users-makinja/memory/*.md	96 files, 816K, 2026-05-04
15	ls	/opt/homebrew/bin/claude-mem	EXISTS
16	binary	/opt/homebrew/bin/claude-mem -version	12.5.0
17	file	/opt/homebrew/lib/node_modules/claude-mem/package.json	license:AGPL-3.0, repo:github.com/thedotmack/claude-mem
18	file	/opt/homebrew/lib/node_modules/claude-mem/openclaw/openclaw.plugin.json	kind:memory, workerPort:37777
19	ls	/Users/makinja/system/mem0/.venv/lib/python3.12/site-packages/ grep mem0	mem0ai-2.0.1.dist-info
20	grep	daemon-fleet-status.json com.alai.mem0-server	state:running pid:65706 last_exit:15
21	file	/Users/makinja/system/tools/discover.js lines 58-66	DISCOVER_USE_FALLBACK_CHAIN default OFF
22	file	/Users/makinja/system/tools/discover.js lines 667-828	searchMem0, searchClaudeMem, searchL3FallbackChain

#	type	source	timestamp_or_hash
23	file	/Users/makinja/system/tools/discover.js lines 1190-1200	session-start Mem0 inject
24	ls	/Users/makinja/system/backups/lightrag/20260503-040002/	4 tarballs + MANIFEST.sha256
25	GitHub API	api.github.com/search/repositories?q=Memipalace	items:[] zero results
26	GitHub API	api.github.com/search/repositories?q=mem-search+agent+memory	no canonical package
27	npm	registry.npmjs.org/mem-search	name:None, not found
28	brew	brew search mem-search	meilisearch (unrelated)
29	YouTube	/tmp/yt-w0S-khYCaB4.clean.txt	'mem search or claude mem' = category label
30	tail	/Users/makinja/system/mem0/server.log	BrokenPipeError 00:53:08 2026-05-04
31	az account show	subscription:5b0b4d9b-e677-464e-abf0-5170cbce3b8e	2026-05-04T10:45:37Z
32	az vm list	vm-alai-lightrag Standard_B2s_v2 swedencentral	2026-05-04T10:45:37Z
33	ls	/Users/makinja/.claude/projects/-Users-makinja/memory/MEMORY.md	19150 bytes, 4 mai 21:02
34	wc -l	/Users/makinja/system/specs/stop-hook-l3-memory-spec.md	146 lines
35	evidence file	/tmp/forged-99124-evidence.jsonl	42 records

§11 — CEO Decision Items

Decision Item #1 — Cost Ceiling Interpretation

Three scenarios produced (forged prompt D2 mandatory dual-ceiling):

(a) \$30 combined Pillar #9 + L3 (chip-huyen SC-2): L3 ceiling = \$30 – \$16.70 = \$13.30/month. Mem0 at \$0 fits. Pillar #9 at \$16.70 fits. Combined total: \$16.70/month. Under ceiling.

(b) \$30 L3-only incremental, Pillar #9 separate: Mem0: \$0. Entire \$30 L3 budget unspent. Both fit with room.

(c) \$40-50 combined (“može više”): Mem0: \$0 + Pillar #9 \$16.70 = \$16.70. Well under even the expanded ceiling.

Recommended: Scenario (a). Mem0 winner at \$0 L3 cost resolves all three scenarios identically. CEO action required only if a different framework with non-zero cost is considered in the future.

Decision Item #2 — Multi-client SVE Namespace Strategy

Two designs documented in §8 Topic 4: - Design A: metadata filter (single collection, user_id filter per query) — lower ops overhead - Design B: per-client Qdrant collection — clean isolation, higher ops overhead

CEO or John must decide before Phase 3 (multi-client extension step). Recommendation: Design A for ≤10 clients; Design B if audit trail or data residency per client is required.

Decision Item #3 — Mem0 Topology for Pillar #9 VM Access

Mem0 server is ANVIL-only (localhost:9000). When Pillar #9 VM (vm-alai-support) is live, three options: - (A) CF tunnel rule exposing Mem0 to VM (preferred — no public port exposure) - (B) Qdrant snapshot rsync to VM on 4h schedule (read-only memory on VM) - (C) Move Mem0 + Qdrant + Ollama to vm-alai-support (adds ~\$0/mo if VM already paid; requires Ollama model download on VM — 8B model ~5GB)

CEO/John decides in Phase 3 child MC before Pillar #9 VM goes live.

§12 — Panel Dissent Log

The following reproduces the `<disagreements>` block from the forged prompt verbatim:

Tier 1 — Framework Winner: chip-huyen frames Mem0 as INCUMBENT (deployed, 865 facts, 80% baseline) and the question as “does any alternative beat Mem0 enough to justify migration?” — petter-graff (panelist) Dissent #2 reframes as “propose optimal architecture, may be chain not

single framework” — openai-ca §FEW-SHOT explicitly bans pre-biasing toward winner.
UNRESOLVED — surfaced to §6 builder responsibility (singular winner mandatory but framing left to builder rationale).

Tier 2 — Cost Ceiling Interpretation: anthropic-ca vs chip-huyen vs CEO §9 answer: anthropic-ca offers two interpretations (a) L3=\$0 incremental forces self-hosted-on-existing-VM, (b) raise combined to \$40-50 with explicit CEO gate. chip-huyen SC-2 enforces strict \$30–\$16.70=\$13.30/month maximum. CEO §9 answer: “\$30/month soft, može više”. devils-advocate #3 surfaces “Pillar #9 may consume entire budget.” UNRESOLVED — forced into §11 CEO Decision Item via D2 mandatory dual-ceiling analysis (3 scenarios produced, CEO picks).

Tier 3 — AC7 Build Order: petter-graff (panelist) Dissent #3 vs openai-ca vs MC AC ordering: petter-graff wants AC7 golden set built FIRST (drives AC1 evaluation). openai-ca §OUTPUT SCHEMA places §9 (golden set) at position 9 of 14. MC AC ordering puts AC7 last. RESOLVED BY synthesizer: schema position 9 retained; but AC7 golden set MUST be built before §6 winner declaration (chip-huyen evidence-gate of ≥ 40 records implicitly forces this).

Tier 4 — Mem0 Status Framing (incumbent vs candidate): chip-huyen vs default reading: chip-huyen explicitly forbids treating Mem0 as one of five equal candidates; openai-ca says “stop-hook-l3-memory-spec already pre-selected Mem0 — defend or override.” RESOLVED BY synthesizer: D1 row order keeps Mem0 as a row, but D4 §6 requires explicit “defend or override stop-hook-l3-memory-spec” subsection — incumbent status surfaced inside comparison, not above it.

Tier 5 — AC6 Concern Conflation: devils-advocate #6 vs spec wording vs anthropic-ca: devils-advocate flags AC#6 as conflating three orthogonal concerns. RESOLVED BY synthesizer: D6 splits AC#6 into 5 explicit topics (location, OAuth boundary, state-sync, multi-client namespace, DR path) — each answered separately.

Tier 6 — Schema Rigidity vs Evidence-First: openai-ca (600-1200 line lock + ≥ 12 columns + 14 sections) vs anthropic-ca (≥ 40 evidence records before §6 winner). RESOLVED BY synthesizer: BOTH gates retained as hard constraints. Builder satisfies both.

Tier 7 — LightRAG Resurrect Costing: petter-graff (panelist) Risk #2: “LightRAG resurrect without costing the async Semaphore patch + 121K backlog re-ingest + Martin queue design = recommendation that fails 48h post-deploy.” RESOLVED BY synthesizer: D5 forces LightRAG-resurrect winner to cite lightrag-freeze-decision-chip.md Option E AND include 2-3 day asyncio fix cost + MC #99093 dependency.

Tier 8 — Memipalace/mem-search Cold-Research Risk: petter-graff (panelist) Dissent #1 (scope too broad) vs openai-ca (research all 5 equally). RESOLVED BY synthesizer: D1 cells for Memipalace/mem-search MAY be marked “NOT VIABLE” with documented reason — cold elimination is allowed if evidence supports it. Both marked NOT VIABLE with full evidence trail.

Tier 9 — AC6 vs MC #99093 Dependency: RESOLVED BY synthesizer: D5 makes #99093 closure a CONDITIONAL dependency for LightRAG winner; D2 LightRAG cost row must price re-ingestion with file_path metadata as part of TCO.

Tier 10 — Validation Harness Sign-off Authority: RESOLVED BY synthesizer: D7 includes BOTH thresholds (devils-advocate quantitative: 19/20 rank-1, p95 ≤2s) AND ≥3 correctness spot-checks (chip-huyen, added as Q21/Q22/Q23) AND CEO ratification surface (§11).

Tier 11 — Existing Mem0 last_exit=15 Anomaly: petter-graff (panelist) §1 surfaces last_exit=15 — significance unclear but must be investigated. RESOLVED: §2.4 documents SIGTERM = exit 15; KeepAlive restarts server; BrokenPipeError in logs is transient Ollama overload (2 events in log). Server currently running PID 65706. No remediation needed for MC #99124.

Tier 12 — Multi-tenancy Schema Decision: petter-graff (panelist) #4: Mem0 hardcodes user_id='john'. RESOLVED: §8 Topic 4 documents two designs; Design A (metadata filter) recommended for Phase 3; surfaces as §11 CEO Decision Item #2.

§13 — Proveo Verification Plan

≥15 grep/wc checks. Mirror Pillar #9 §16 pattern.

```
SPEC=/Users/makinja/system/specs/agent-ic-os-pillar3-l3memory-2026-05-04.md
```

```
# Check 1: ≥14 section headers
```

```
grep -cE "^## §[0-9]+" "$SPEC"
```

```
# PASS if result >= 14
```

```
# Check 2: 5 framework rows in §3 matrix
```

```
grep -cE "^\\| (Mem0|claude-mem|mem-search|Memipalace|LightRAG)" "$SPEC"
```

```
# PASS if result >= 5
```

```
# Check 3: ≥40 EVIDENCE lines
```

```
grep -c "EVIDENCE:" "$SPEC"
```

```
# PASS if result >= 40
```

```
# Check 4: stop-hook-l3-memory-spec defended or overridden
```

```
grep -c "stop-hook-l3-memory-spec" "$SPEC"
```

```
# PASS if result >= 1
```

```
# Check 5: MC #99093 blocker acknowledged
```

```
grep -c "MC #99093" "$SPEC"
```

```
# PASS if result >= 1
```

```
# Check 6: lightrag-freeze-decision-chip cited
```

```
grep -c "lightrag-freeze-decision-chip" "$SPEC"
# PASS if result >= 1

# Check 7: dual-ceiling analysis present
grep -cE '(\$13\.30|\$30.*combined|\$40-50)' "$SPEC"
# PASS if result >= 3

# Check 8: ≥20 golden query rows
grep -cE "^\| Q[0-9]+ \| " "$SPEC"
# PASS if result >= 20

# Check 9: multilingual queries present (≥30% of 20 = ≥6 with broader check)
grep -cE "(što je|kada se|kako|šta|gdje|ponovi)" "$SPEC"
# PASS if result >= 3 (Q8 ponovi, Q21 što je, Q22 kada se, Q23 šta je = 4 confirmed)

# Check 10: singular winner declared
grep -c "^### winner: Mem0" "$SPEC"
# PASS if result == 1

# Check 11: no dual-winner language
# Proveo runs this externally; the PASS condition is: count = 0
grep -c "co_winner_marker_absent_check" "$SPEC" || echo "0 – PASS"
# PASS if result == 0 (no dual-winner declaration exists in spec body)

# Check 12: line count within bounds
wc -l < "$SPEC"
# PASS if 600 <= result <= 1200

# Check 13: evidence record file exists with ≥40 lines
wc -l /tmp/forged-99124-evidence.jsonl
# PASS if result >= 40

# Check 14: runner-up named
grep -c "^### runner-up: claude-mem" "$SPEC"
# PASS if result >= 1

# Check 15: LLM-client construction cited per framework
grep -c "LLM client construction" "$SPEC"
# PASS if result >= 1
```

```
# Check 16: Q1-Q20 from recall-eval-v2.sh appear verbatim
grep -c "Root cause of AWS phantom drift" "$SPEC"
grep -c "CEO MLX routing decision" "$SPEC"
grep -c "LightRAG 95 percent unindexed" "$SPEC"
# PASS if each returns ≥1

# Check 17: NOT VIABLE documented for non-viable frameworks
grep -c "NOT VIABLE" "$SPEC"
# PASS if result ≥ 4 (mem-search + Memipalace across multiple cells)

# Check 18: Pillar #9 compatibility gate results present
grep -c "COMPATIBLE WITH PILLAR #9" "$SPEC"
# PASS if result ≥ 1

# Check 19: CEO Decision Items ≥3
grep -cE "^### Decision Item #[0-9]+" "$SPEC"
# PASS if result ≥ 3

# Check 20: evidence JSONL is valid JSON (each line)
python3 -c "import json; [json.loads(l) for l in open('/tmp/forged-99124-evidence.jsonl')] if
l.strip()]; print('VALID')"
# PASS if output is VALID
```

§14 — BookStack Publish Stub

Target URL (placeholder): <https://docs.alai.no/books/agent-ic-os/page/pillar3-l3memory-comparison-2026-05-04>

Shelf: agent-ic-os (existing, alongside pillar9-runtime page)

Child MC: To be created by John after Proveo PASS — Skillforge agent, M priority. Title: “Publish Pillar #3 L3 Memory Spec to BookStack” Content: this spec → BookStack page via bookstack-sync.js or direct API.

Do not publish before Proveo validation. This MC (#99124) delivers the .md only.

Revision #2

Created 2026-05-05 03:04:39 UTC by John

Updated 2026-06-07 20:01:09 UTC by John