

# ALAI AI System — v2.0 Operating Picture & Master Roadmap

# ALAI AI System — v2.0 Operating Picture & Master Roadmap

**Date:** 2026-05-19 **Architect:** Petter Graff **Status:** SYNTHESIS COMPLETE — pending dual validation (Proveo + Verifier) **Supersedes:** `ceo-ai-system-audit-2026-05-18-REPORT.md` (v1.1 — Wave 1 still canonical for inventory; v2.0 adds design + build roadmap)

---

## 1. Executive Brief

The ALAI AI system is a *system that builds systems* — and it has stopped building. Over the last 8 days it burned **\$742K on Anthropic Opus (99.98% of all spend)**, peaked at **\$377,487 in a single day (2026-05-11)**, and shipped **zero production code in 7 days**. Wave 1 (2026-05-18) identified the symptoms; Wave 2 (three parallel teams: Control, Knowledge, Workflow) identified the single causal narrative:

“ **The orchestrator steers by frozen instruments, dispatches through gates that don't fire, into a free-tier fleet that doesn't exist, validates with probes that never run, and ships into a backlog with no exit.** Every "save" is a watchdog that itself is dormant. The meta-failure — `hook-drift-detector` daemon exit 2, stopped — is what allows all other silent failures to hide.

The three planes fail compoundingly:

- **Control plane:** `opus-cost-guard` has no daily \$ ceiling, defaults ALLOW when `model` field is absent, doesn't gate the main session — only sub-Tasks. The May 11 \$377K spike *would not have been blocked*. 4 of 14 tier-routes are ghosts (devstral:24b absent, 2/3 MLX serve wrong model = bge-m3). Most hooks have zero audit logs today (verifier: 60 hooks on disk, majority dark). Evidence ledger SQLite has 0 tables; the JSONL has 107 verdict rows, 79/107 (74%) `force_completion` and 0 `PROBE_PASS` — gate-gaming theater (verifier-corrected).
- **Knowledge plane:** Mem0 (Pillar #3 winner per project\_99124) is dead in runtime (port 9000=000, no LaunchAgent). `discover.js` cites `manifest-index.md` (mtime 2026-04-06, 43 days stale; embedded audit date 2026-02-26). `skill-registry.db` carries 96 skill rows but only 12 with non-zero `use_count` and no `last_used` column. BookStack API blocked (CF Access 302). LightRAG pump hard-capped at 600/run with 23,558 backlog that grows. ZAKON #12 RAG injection is referenced but unwired — every dispatch re-inhales ~15K-token MEMORY.md.
- **Workflow plane:** 873 of 887 emails (98.4%) unlinked to MC tasks. `discover.js routing` CLI cited in CLAUDE.md **does not exist** — routing is improvised by LLM. `mehanic` + `dzevad-jahic` referenced but absent from `specialist-mapping.json`. `claude-builder durable-runner`: 2,945 failed / 1 completed since April. 2,400 zombie MC tasks >14d. TLDR daemon writes to `~/system/data/insights/` which does not exist.

## If you read nothing else

- **A single \$-ceiling hook (T-A-02) ships in 1 day and would have prevented the entire May 11 spike. Build it first.**
- **The control plane must turn on before the knowledge plane gets fixed before the workflow plane closes the loop. Week 1 → Week 2 → Week 3.**
- **9 CEO decisions are surfaced (\$6). Six are go/no-go on existing components; three are scope-of-resumption.**
- **Conservative combined save: \$780K-\$2.7M/month. Build cost: <\$100. Payback <1 hour of current burn.**

## One sentence per plane

- **Control:** Today blind & ungated → Week 1 kill-switch + \$-ceiling + tier reconcile + Reality Anchor watchdog.
  - **Knowledge:** Today stale & lying → Week 2 CF token + ZAKON #12 wire + manifest regen + 8 governance pages on BookStack.
  - **Workflow:** Today disconnected end-to-end → Week 3 email→MC daemon + `router.js` + TLDR + backlog TTL + escalation matrix.
  - **Production code:** Resumes Week 4 only after E2E test (CEO email → done in <90 min, no mid-loop prompts) passes 8/9.
-

# 2. The Three Planes (Target Architecture)

## 2.1 Mermaid Super-Diagram

```
flowchart TB
  subgraph CEO_SURFACE [CEO Surface]
    Prompt[CEO prompt / Slack]
    Email[CEO email IMAP]
  end

  subgraph CONTROL [Plane 1 – Control & Determinism]
    KS[Kill switch<br/>tmp alai-killswitch]:::new
    OCG[opus-cost-guard v2<br/>daily $ ceiling]:::fix
    KSW[fleet-reconcile-probe<br/>tier-truth.json]:::new
    RAW[probe-liveness-watchdog]:::new
    HDD[hook-drift-detector v2]:::new
    EL[(evidence-ledger.db<br/>SQLite schema'd)]:::fix
    SSM[session-spend-monitor<br/>per-session $ ladder]:::new
  end

  subgraph KNOWLEDGE [Plane 2 – Knowledge & Memory]
    DJ[discover.js<br/>3-tier front door]:::fix
    L1[L1 MEMORY.md + session]:::ok
    L2[L2 HiveMind 21,741 rows]:::ok
    L3a[L3a LightRAG Azure]:::fix
    L3b[L3b Mem0 facts<br/>KILL → fold to HiveMind]:::kill
    BS[(BookStack 478 pages<br/>canonical wiki)]:::fix
    Z12[ZAKON #12<br/>rag-context-for-builder]:::new
    INV[manifest-index + skill-registry<br/>daily regen]:::fix
  end

  subgraph WORKFLOW [Plane 3 – Orchestration & Workflow]
    EID[email-intake-daemon]:::new
    MC[(MC tasks db)]:::ok
    RTR[router.js classify<br/>discover.js routing alias]:::new
    MEH[mehanik gate]:::fix
  end
```

```
SUB[Specialist subagents]:::ok
PIO[pi-orchestrator<br/>route_eligibility expanded]:::fix
PRO[Proveo E2E validation]:::ok
TLDR[TLDR daemon<br/>~/system/data/insights]:::new
TTL[backlog-ttl-daemon]:::new
ESC[escalation-matrix hook]:::new
end
```

```
Prompt --> Z12
Email --> EID --> MC
MC --> RTR --> MEH --> SUB
SUB -.queries.-> DJ
DJ --> L1 & L2 & L3a & L3b
DJ -. cite .-> BS
Z12 --> DJ
SUB --> OCG
OCG -. breach .-> KS
SSM -. breach .-> KS
KS -. blocks.-> SUB & MEH
KSW -. health .-> SUB
RAW -. probes .-> PRO
PRO --> EL
EL --> MC
HDD -. watches .-> OCG & KSW & RAW & EID & TLDR
PIO --> PRO
SUB --> PIO
MC --> TTL
TTL --> TLDR --> Prompt
ESC -. gates .-> Prompt
INV -. truth .-> DJ
```

```
classDef new fill:#1d8c43,color:#fff
classDef fix fill:#d4a017,color:#000
classDef kill fill:#b3261e,color:#fff
classDef ok fill:#5b9bd5,color:#fff
```

Legend: green = new build, yellow = fix-in-place, red = formal kill, blue = working today.

## 2.2 Plane Summaries

**Control plane (Team A).** *Current:* Probes designed but not running (0 PROBE\_PASS events 7d). Hooks present (58) but only 5 with today's audit logs. `opus-cost-guard` blocks per-agent name match, not \$-ceiling. May 11 (\$377K) would not have triggered any gate. Evidence ledger SQLite empty (0 tables); JSONL = 100% `force_completion`. Tier router blind: 4/14 routes point at ghost models. *Target:* Hard \$-ceiling + global kill-switch + live fleet reconcile (5-min cycle) + Reality Anchor watchdog auto-restarting dormant probes + evidence-ledger schema with HMAC chain + per-hook audit-log convention enforced by hook-drift-detector v2. *MCs:* 9 (T-A-01 through T-A-09).

**Knowledge plane (Team B).** *Current:* 5 critical governance subsystems (Reality Anchor, ZAKON NULA, Tier Router, Evidence Ledger, Hooks) have ZERO BookStack pages. `discover.js` cites stale manifest. ZAKON #12 dormant — every builder dispatch eats ~15K tokens of full MEMORY.md re-injection. LightRAG: degraded (15% timeout), public endpoint CF Access blocked, pump capped 600/run with 23,558 backlog. Mem0 dead. ADR numbering collisions (025×2, 026×4). *Target:* One front door (`discover.js memory --budget=2000`) that spans L1+L2+L3 with token-budget contract. CF Access rotated → BookStack + LightRAG public both unblocked. ZAKON #12 wired into PreToolUse → ~105K tokens/day saved. 8 governance pages published; ADR allocator + collision repair. Mem0 killed (Path B), folded into HiveMind facts table. Library built (Path A) as central skill registry. *MCs:* 17 (MC-B01 through MC-B17).

**Workflow plane (Team C).** *Current:* CEO email pipeline broken at every transition. Email→MC linkage dead (873/887 unlinked, 80 `replay_required` with no replay daemon). `discover.js routing` CLI is fictional. claude-builder queue: 2,945 failed since April. PI-orch alive but `route_eligibility=['post-build']` excludes every real MC. TLDR daemon writes to nonexistent dir. 2,400 zombie MCs. 65 agent files vs 30 mapping keys. *Target:* `email-intake-daemon` classifies via local qwen3 (\$) → MC link 100%. `router.js classify` made real (alias makes CLAUDE.md claim honest). Mapping JSON closed (0 orphans). `backlog-ttl-daemon` enforces 30d/60d retirement. PI-orch route filter expanded to 5 categories → free-tier execution path revived. Session-spend-monitor closes the gap opus-cost-guard cannot (main session burn). Escalation matrix hook silences micro-decision pings to CEO. *MCs:* 13 (MC-C1-1 through MC-C5-1).

---

## 3. Cross-Plane Couplings (the new picture Wave 1 didn't see)

These five couplings are why no single team can finish in isolation, and why sequencing matters.

### 3.1 ZAKON #12 wire-in = A + B + C all three

- **A owns** the PreToolUse hook plumbing (`~/claude/settings.json` registration, audit log convention from T-A-08). Source: `team-a/control-plane-build-plan.md` T-A-08 + cross-team note line 182-184.

- **B owns** the retrieval logic — `rag-context-for-builder.js` rewrite with `--tier-budget L1:1200,L2:500,L3:300 --max-tokens 2000` (MC-B04). Source: `team-b/knowledge-plane-design.md` §3 + `team-b/knowledge-plane-build-plan.md` MC-B04/MC-B05.
- **C consumes** — every specialist dispatch through the new pipeline receives the 1,800-token block instead of MEMORY.md (workflow plane §3 sequence diagram). Source: `team-c/workflow-plane-design.md` §3.
- **Coupling rule:** B's MC-B05 cannot ship until A's hook framework lands; C's MC-C1-2 router classification reads the same `specialist-mapping.json` that B's MC-B16 patches. **Sequence: A finishes hook framework day 7 of Week 1 → B ships MC-B04/B05 Week 2 → C dispatches through both Week 3.**

## 3.2 Cost guard is 3 layers, one per plane

- **A — gate:** `opus-cost-guard v2` PreToolUse[Task] hard-block on daily \$ ceiling + flip ALLOW-on-missing-model default to BLOCK. Source: `team-a/control-plane-design.md` COMP-1 + `team-a/control-plane-audit.md` §3 "CRITICAL GAP 1-4".
- **B — token-budget:** `rag-context-for-builder --max-tokens` ceiling per dispatch (105K tokens/day saved). Source: `team-b/knowledge-plane-design.md` §3 "Token-save math".
- **C — session ceiling:** `session-spend-monitor.js` polls `costs.db` by `session_id` every 5 min, Slack at \$200 / model-flip at \$500 / kill at \$1,000. This **closes the gap A cannot reach** because `opus-cost-guard` fires on Task subagent dispatch but not on the main session. Source: `team-c/workflow-plane-audit.md` §9 + `team-c/workflow-plane-design.md` §2.5 + `team-c/workflow-plane-build-plan.md` MC-C2-2.
- **Coupling rule:** All three must land. A alone leaves the main session burning; B alone leaves the gate-bypass open; C alone has no per-dispatch ceiling.

## 3.3 `discover.js` is the single front door — three teams patch it

- **A doesn't touch** `discover.js` directly but its T-A-03 `tier-truth.json` becomes a tier health source for B's L3 latency budgeting.
- **B regenerates** `manifest-index.md` + `skill-registry.db` daily (MC-B06), adds `--self-check` meta-probe at boot (MC-B07), upgrades `discover.js memory` to span 3 tiers (MC-B08). Source: `team-b/knowledge-plane-design.md` §7.
- **C makes** `discover.js routing` claim true via `router.js classify` alias (MC-C1-2). Source: `team-c/workflow-plane-audit.md` Break #2 + `team-c/workflow-plane-design.md` §2.2.
- **Coupling rule:** John currently does tool-first verification through a `discover.js` that lies; until all three patches land (B inventory regen + C routing alias), every "tool-verified" claim downstream inherits residual rot.

## 3.4 Email pipeline is ONE workflow with THREE breaks

The CEO daily flow has a single physical pipeline (Email → email-inbox.db → MC → router → mehanik → specialist → proveo → done → TLDR) with three independent breaks:

- (B→E) Email-to-MC linkage broken (873/887 unlinked) — team-c/workflow-plane-audit.md Break #1.
- (F) discover.js routing CLI fictional — Break #2.
- (J) TLDR daemon writes to nonexistent ~/system/data/insights/ — Break #4.
- **Coupling rule:** Fixing only one keeps the pipe dark. **MC-C1-1 + MC-C1-2 + MC-C1-4 must ship as a triple** in Week 3 days 1-3. Without all three, CEO email "Pls fix Bilko 500" never reaches a specialist.

## 3.5 Gate-gaming (verdict-ledger 100%

force\_completion) is a consequence of A + B + C all failing

- **A** — probes off → no PROBE\_PASS rows → only path to "done" is --force. Source: team-a/control-plane-audit.md §5 "107 rows, all force\_completion".
- **B** — discover.js lies → builder doesn't know correct evidence path → fabricates artifact (Proveo hallucination 2026-05-07). Source: MEMORY.md feedback\_proveo\_hallucination\_2026-05-07.md.
- **C** — claude-builder queue dead → fallback to inline subagent → no durable record → trivial to fake claim. Source: team-c/workflow-plane-audit.md Break #5.
- **Coupling rule:** "Stop gate-gaming" is **not a single-MC fix**. The fix is sequential: T-A-06 Reality Anchor watchdog → T-A-07 evidence ledger schema + null-path block at mc.js done → MC-B04 ZAKON #12 wire (so builders get correct context) → MC-C1-1 email→MC (so MCs land with real source) → MC-C4-2 claude-builder fossil archive. After this chain, verdict-ledger PROBE\_PASS:force\_completion ratio shifts from 0:107 toward 50:50 within 7 days (T-A-06 AC).

## Cross-Team Contradictions (resolved)

Reviewed all three audit docs for conflicting claims; **no hard contradictions found**, only resolved revisions:

- **Team C corrects Wave 1 on PI-orch.** Wave 1 said "pi-orch HTTP dead 50d"; Team C probed launchctl list and found PID 57544 alive, polling, but route\_eligibility=['post-build'] matches zero real MCs. **Verdict:** PI-orch is alive but useless; the underlying claim ("free-tier execution path is broken") holds. Memory note project\_ai\_factory\_audit\_2026-05-09 should be updated.
- **Team C corrects Wave 1 on skill-registry.** Wave 1 said 1 row; Team C found 96 rows (registry was rebuilt at some point) but only 12 have non-zero use\_count and there's no last\_used timestamp — so the substantive claim ("skill catalog isn't measured") holds.

- **Team C corrects Wave 1 on edita queue.** Wave 1 cited 161 dead-letter; Team C found 22 in `dead_letter_queue` but 2,945 in `queue_entries` failed against `claude-builder`. The number moved tables; the magnitude is **larger**, not smaller.

## 4. Master Roadmap (4 Weeks)

Week	Theme	Teams	MCs to ship	End-state gate (deterministic probe)	Rollback
1	Stop the bleed	A	T-A-01 kill switch, T-A-02 \$ ceiling, T-A-03 fleet reconcile, T-A-04 devstral, T-A-05 MLX, T-A-06 probe watchdog, T-A-07 evidence schema, T-A-08 hook-drift v2, T-A-09 daemon sweep	<code>control-plane-health.sh</code> returns 7/7 PASS: killswitch round-trip; cost-ceiling fires at synthetic \$1000; tier-truth.json all 14 tiers healthy or explicitly disabled; probe-watchdog detects 48h synthetic stall; evidence-ledger.db has table + row-count == JSONL; hook-drift detects 24h synthetic silence; 0 flapping daemons	Disable killswitch + revert hook-drift v2 plist; T-A-02 ceiling can be raised to \$10K/day as soft-rollback. Evidence schema is additive — no rollback needed.
2	Lights on	B (+ A finishing T-A-08 integration)	MC-B01 CF token, MC-B02 LightRAG pump, MC-B03 outbox-ingest decision, MC-B04 rag-context rewrite, MC-B05 ZAKON #12 wire, MC-B06 inventory regen, MC-B07 self-check, MC-B08 memory upgrade, MC-B09 HiveMind purge, MC-B10 dead-agent TTL	<code>discover.js --self-check</code> reports 0 drift on day 7; <code>curl https://lightrag.alai.no/health</code> returns 200; <code>bookstack-staleness.js sample</code> returns JSON; ZAKON #12 fires logged for ≥80% of builder dispatches; pre/post token count shows ≥40% reduction in builder prompts	MC-B05 hook is opt-in via env flag <code>ZAKON12_ENABLED=1</code> for first 24h; if drift >5% on day 1, revert to off. MC-B09 stub removal: archive-first, restore is <code>cp</code> from <code>_archive/</code> .

Week	Theme	Teams	MCs to ship	End-state gate (deterministic probe)	Rollback
3	Workflow restored	C	MC-C1-1 email→MC, MC-C1-2 router.js, MC-C1-3 mapping cleanup, MC-C1-4 TLDR, MC-C2-1 backlog TTL, MC-C2-2 session-spend, MC-C2-3 per-MC budget, MC-C3-1 HiveMind cleanup, MC-C3-2 skill registry, MC-C3-3 MCP cleanup, MC-C4-1 pi-orch routes, MC-C4-2 claude-builder archive, MC-C5-1 escalation hook	<b>E2E test:</b> CEO sends 1 test email → MC linked <5min → routed → mehanik authorized → specialist returned <60min → Proveo PASS to Slack #ceo-digest with screenshot → TLDR digest 6h later. 8/9 sub-criteria pass.	MC-C1-1 daemon can be disabled; backfill MC link via one-off script. MC-C2-2 session monitor is alert-only first 48h before model-flip is enabled. MC-C5-1 hook is WARN-only first 7 days.
4	Production resumes	All teams hardening + Bilko/Drop work	Production MCs from BUILD-BLUEPRINT.md per project; no new system-level MCs except hardening	<pre>git log --since=7.days --author=alai-builders ~/projects/bilko-cloud &gt; 5 commits AND costs.db today &lt; \$5K AND verdict-ledger PROBE_PASS:force completion ≥ 1:1</pre>	If Week 4 cost burn returns to >\$10K/day → freeze prod work, return to Week 3 hardening. Killswitch always available.

**Gate between weeks:** each week's end-state probe must PASS before the next week's specialist dispatches are authorized. CEO sign-off on probe report = go.

## 5. MC Inventory (Consolidated 39 MCs)

ID	Title	Team	Prio	Week	\$ Save	Dep
T-A-01	Kill switch + CLI	A	BLOCKER	1	insurance	—
T-A-02	opus-cost-guard v2 daily \$ ceiling	A	BLOCKER	1	\$20-70K/d	T-A-01

ID	Title	Team	Prio	Week	\$ Save	Dep
T-A-03	fleet-reconcile-probe + tier-truth	A	H	1	\$2-8K/d	T-A-01
T-A-04	devstral pull or remap	A	H	1	\$5-15K/d	T-A-03
T-A-05	MLX M2c+M3 repair	A	H	1	\$1-5K/d	T-A-03
T-A-06	Reality Anchor watchdog	A	H	1	risk-redux	T-A-01
T-A-07	Evidence ledger SQLite schema	A	H	1	risk-redux	—
T-A-08	hook-drift-detector v2	A	M	1	risk-redux	T-A-01, T-A-07
T-A-09	Daemon hygiene sweep	A	M	1	\$0 direct	—
MC-B01	CF Access token rotate	B	H	2	unblock \$15-42/mo	—
MC-B02	LightRAG pump 600→5000	B	H	2	40-80K tok/d	B01
MC-B03	outbox-ingest restore/decom (ADR-036)	B	M	2	qual	B01
MC-B04	rag-context-for-builder rewrite	B	H	2	105K tok/d	B02, T-A-08
MC-B05	ZAKON #12 PreToolUse hook	B	H	2	activates B04	B04, T-A hook fw
MC-B06	Daily inventory regen cron	B	H	2	5-30K tok/d	—
MC-B07	discover.js --self-check at boot	B	H	2	indirect	B06
MC-B08	discover.js memory 3-tier upgrade	B	M	2	qual	B02, B06

ID	Title	Team	Prio	Week	\$ Save	Dep
MC-B09	Purge 3 orphan HiveMind stubs	B	M	2	10K tok/d	—
MC-B10	Dead-agent TTL ADR-035	B	M	2	6K tok/d	—
MC-B11	bookstack-staleness daemon revive	B	H	3	\$0 direct	B01
MC-B12	Publish 8 governance pages	B	H	3	\$0 direct	B01
MC-B13	ADR allocator + 6 collision repair	B	M	3	\$0	—
MC-B14	Mem0 ADR-033 (recommend KILL)	B	M	3	consolidation	—
MC-B15	Library ADR-034 (recommend BUILD)	B	M	3	qual	B06
MC-B16	specialist-mapping audit	B	M	3	\$1-3/mo	B06
MC-B17	Hook .bak cruft cleanup	B	L	3	\$0	—
MC-C1-1	email-intake-daemon	C	BLOCKER	3	unblock A	T-A fleet
MC-C1-2	router.js classify CLI	C	H	3	unblock	C1-3
MC-C1-3	specialist-mapping completion + ADR-027	C	H	3	\$1-3/mo	—
MC-C1-4	TLDR daemon reconnect	C	H	3	qual (closes loop)	C1-1
MC-C2-1	backlog-ttl-daemon	C	H	3	signal/noise	C1-4
MC-C2-2	Session spend monitor (Layer 2)	C	BLOCKER	3	\$5-30K/d session cap	T-A-02

ID	Title	Team	Prio	Week	\$ Save	Dep
MC-C2-3	Per-MC budget (Layer 3)	C	H	3	\$1-5K/d	C2-2
MC-C3-1	HiveMind ~85 zombie + 46 pollution cleanup	C	M	3	qual	—
MC-C3-2	Skill registry + retire wave	C	M	3	qual	—
MC-C3-3	MCP audit + decom stitch+local-rag (ADR-029)	C	M	3	startup time	—
MC-C4-1	pi-orch route_eligibility expansion	C	M	3	free-tier revival	T-A-04, T-A-05
MC-C4-2	claude-builder fossil archive (ADR-030)	C	M	3	\$0	—
MC-C4-3	edita owner audit + reassign	C	M	3	signal/noise	—
MC-C5-1	Escalation matrix hook	C	H	3	CEO-attention save	C1-4

Plus 5 Wave 1 P0 carryovers (now subsumed): P0-1 #101375 → T-A-02; P0-2 #101376 → T-A-04; P0-3 #101377 → T-A-06; P0-4 #101378 → MC-B07; P0-5 #101379 → T-A-05.

**Total Wave 2 MCs:** 40 distinct (including MC-C4-3) + 5 Wave 1 P0 consolidated.

## 6. Risks & Open CEO Decisions

- Mem0 — resurrect (Path A) or kill+fold-into-HiveMind (Path B)?** *Recommendation: B.* Reduces moving parts; Qdrant runtime removed; HiveMind `facts` table covers same use case. Mem0 has been dead 14+ days with no detected loss. Formalize via ADR-033 (MC-B14).
- Library system — build (Path A) or kill (Path B)?** *Recommendation: A — minimal build.* `~/system/library.yaml` is real intent, no consumer ever shipped. A 1-day install script gives one-place control over which skills are active where; the alternative is 96 skills with no source-of-truth. Formalize via ADR-034 (MC-B15).
- PI-orchestrator — expand route filter (Path A) or formal decommission (Path B)?** *Recommendation: A first, B as fallback.* MC-C4-1 expands `route_eligibility` to 5 categories. **Kill criterion (auto):** if after T-A-04 + T-A-05 + MC-C4-1 ship, pi-orch still has

0 matching tasks in 7 days, formal kill via ADR-026 (one of the existing collision files — repaired in MC-B13).

4. **claude-builder durable-runner queue — drain + restart, or replace?**

*Recommendation: drop the queue, do not restart.* 2,945 failed / 1 completed since April = the architecture is fossilized. MC-C4-2 archives. Future "durable-runner v2" decision punts to Week 5+; not in current scope.

5. **2,400 zombie MC tasks — auto-close at >14d idle?** *Recommendation: tiered TTL via MC-C2-1.* Open + M/L + >30d → auto-pause. Paused + >60d → auto-close. H + open + >14d → CEO digest entry. **Not** blanket auto-close — preserves CEO-owned tasks (alem has 72 open).

6. **Production code resumption — Week 4 firm or conditional?** *Recommendation: conditional on Week 3 end-state E2E probe (8/9 sub-criteria PASS + 48h cost <\$5K/day).* If both gates green, resume Week 4. If either red, Week 4 = hardening cycle; production code Week 5.

7. **Daily \$ ceiling level (T-A-02) — \$500/day Opus default?** *Recommendation: yes, with `~/system/config/cost-ceilings.json` knob.* Pre-AI-Services-revenue, \$500/day Opus = \$15K/month. Override token TTL 60s for CEO-explicit cases. If CEO wants \$300/day, change one JSON line.

8. **Session-spend ladder (MC-C2-2) — \$200 alert / \$500 model-flip / \$1000 kill?** *Recommendation: alert-only first 48h, then enable model-flip + kill.* Avoids same-day surprise on already-running session.

9. **Wave 2 build budget — what's the Opus ceiling for the build phase itself?** *Recommendation: \$250 total for all 40 MCs.* Each MC ≈ \$1 prompt-forge + \$2-5 specialist + \$1 Sonnet sub + \$1 Proveo + \$0.50 Skillforge ≈ \$5-8 avg. Build cost << 1 hour of current burn. Use `/prompt-forge` only for H/BLOCKER (Week 1 + Week 3 BLOCKERS); skip for M/L.

## 7. Total Economics

Source	Daily save (conservative)	Daily save (optimistic)	Monthly (conservative)
T-A-02 cost ceiling	\$20,000	\$70,000	\$600,000
T-A-03/T-A-04 ghost tier kill	\$5,000	\$15,000	\$150,000
T-A-05 MLX repair	\$1,000	\$5,000	\$30,000
MC-B04/B05 ZAKON #12 wire	\$0.50 (token)	\$1.40 (token)	\$15-42 (token equiv)
MC-B06 inventory regen (re-dispatch prevent)	\$0.30	\$1.80	\$9-54
MC-C2-2 session spend ladder (caps catastrophic)	\$5,000	\$30,000	\$150,000

Source	Daily save (conservative)	Daily save (optimistic)	Monthly (conservative)
MC-C1-1 email→MC (operational efficiency)	\$0 direct	\$0 direct	unblocks revenue
MC-C2-1 backlog TTL (signal/noise)	\$0 direct	\$0 direct	CEO time
<b>Total</b>	<b>~\$26,000/day</b>	<b>~\$90,000/day</b>	<b>\$780K-\$2.7M/month</b>

**Wave 2 build phase cost (Opus + Sonnet):** ~\$250 one-time (see Decision 9).

**Payback:** <1 hour of current burn at conservative \$26K/day = \$1,083/hour. Build pays for itself in roughly 13 minutes of current operations.

## 8. Validation Plan

### 8.1 Proveo (Angie Jones) — re-probe ?20% of synthesis claims

Focus areas (load-bearing claims):

- Cross-plane coupling 3.1: ZAKON #12 token-save math (10 dispatches × 10,500 tok). Verify `wc -l` on actual MEMORY.md + measured builder prompt sizes.
- Coupling 3.2: that `opus-cost-guard` does NOT gate main session — re-run probe `~/cache/opus-cost-guard-*.log` for last 48h on current Opus session.
- Coupling 3.4: re-run `sqlite3 email-inbox.db "SELECT COUNT(*) FROM emails WHERE status='new' AND mc_task_id IS NULL"` — assert ≥870.
- Coupling 3.5: verdict-ledger `force_completion` count — assert ≥100, PROBE\_PASS = 0.
- Master roadmap Week 1 gate: probe `~/system/tools/control-plane-health.sh` (does not exist yet — flag if T-A-09 doesn't ship one).
- Decision 4 evidence: re-probe `claude-builder` queue counts — assert ≥2,900 failed and ≤2 completed.

Output: `~/tmp/proveo-v2-operating-picture-validation.jsonl`.

### 8.2 Verifier — atomic-claim decomposition

Decompose into atomic claims:

- All headline facts in §1 Executive Brief.
- Each row of MC inventory table — task ID, team, priority, week, dep correctness.

- Each "\$ save" figure — does it come from a team build plan, and does the math add up?
- Each "Path X recommended" — is there a cited reason in the corresponding team design?

Verdicts per claim: CONFIRMED / PARTIAL / HALLUCINATION. Cost <\$0.50.

## 8.3 Publish

After dual validation PASS → BookStack page "**System Architecture**" book, page "**ALAI AI System v2.0 — Operating Picture & Master Roadmap (CEO Rebuild Brief)**". This becomes canonical; v1.1 (Wave 1) demoted to historical reference.

# 9. Build Phase Dispatch Order (Week 1 only)

Weeks 2-4 dispatch after Week 1 closes (gate from §4).

```

Day 1 (0-4h): /prompt-forge T-A-01 → /mehanik → FlowForge dispatch (Kelsey)
              AC probe: killswitch round-trip + 17 PreToolUse hooks updated.

Day 1 (4-10h): /prompt-forge T-A-02 → /mehanik → FlowForge + Securion review dispatch
               AC probe: synthetic $1,000 cost row → next Opus dispatch BLOCKED + killswitch
               touched.

Day 2:        /prompt-forge T-A-03 → /mehanik → AgentForge + FlowForge dispatch (Georgi +
               Kelsey)
               AC probe: stop ANVIL Ollama → tier-truth marks 3 tiers unhealthy in 5min →
               restart recovers.

Day 3 (parallel A): /mehanik T-A-04 → AgentForge (Georgi) – devstral pull/remap.
Day 3 (parallel B): /mehanik T-A-05 → AgentForge (Georgi) – MLX M2c+M3 repair.
               Skip /prompt-forge for both (M-priority).

Day 4-5:      /prompt-forge T-A-06 → /mehanik → FlowForge + AgentForge dispatch
               AC probe: touch probe last.jsonl mtime=48h → watchdog STALL + restart in 5min.

Day 5-6:      /mehanik T-A-07 → CodeCraft (Bruce Momjian) dispatch (M-priority, no prompt-
               forge).
               AC probe: insert null-path row → mc.js done exits 2 "evidence_path required".

```

Day 6-7: /mehanik T-A-08 → FlowForge + Securion dispatch.  
AC probe: kill pilot-discover-inject.py 24h → drift detector flags in 15min.

Day 7: /mehanik T-A-09 → FlowForge dispatch (daemon sweep).  
Then run `control-plane-health.sh` master probe.  
7/7 PASS → CEO go-ahead for Week 2 Team B dispatch.  
<7 PASS → Week 1 extends by 1-2 days; do NOT proceed to Week 2.

After every dispatch: `/task-postflight` + verifier subagent in bg (per `feedback_active_verifier_pattern_2026-05-14`).

Each MC closes with `mc.js done <id>` only after Proveo PASS + Skillforge BookStack page (ZAKON PLAN).

## END v2.0 OPERATING PICTURE.

### Sources:

- `/tmp/srz-rebuild-2026-05-19/team-a/{control-plane-audit, control-plane-design, control-plane-build-plan}.md`
- `/tmp/srz-rebuild-2026-05-19/team-b/{knowledge-plane-audit, knowledge-plane-design, knowledge-plane-build-plan}.md`
- `/tmp/srz-rebuild-2026-05-19/team-c/{workflow-plane-audit, workflow-plane-design, workflow-plane-build-plan}.md`
- `~/system/specs/ceo-ai-system-audit-2026-05-18-REPORT.md` (v1.1)
- `~/system/specs/srz-rebuild-3-teams-2026-05-19-plan.md` (charter)

# 10. Validation Patches v2 (applied 2026-05-19 after Proveo + Verifier)

**Sources:** `/tmp/srz-rebuild-2026-05-19/proveo-v2-verdict.json`, `/tmp/srz-rebuild-2026-05-19/verifier-v2-report.json`

Patch	Original	Corrected	Source
V2-P1	"skill-registry.db has 1 row for 96 skills"	96 rows, but only 12 with <code>use_count&gt;0</code> ; needs <code>last_used</code> column	verifier KP4
V2-P2	"Build cost: <\$100"	~\$250 (40 MCs × \$5–8 avg, consistent with §6 Decision 9 math)	verifier D4

Patch	Original	Corrected	Source
V2-P3	"8 governance pages on BookStack"	5 governance pages (Reality Anchor, Determinism, Tier Router, Evidence Ledger, Hooks)	verifier KP11
V2-P4	"Total Wave 2 MCs: 39 distinct"	40 distinct (MC-C4-3 edita owner audit was missed in count)	verifier MC1
V2-P5	"65 agent files vs 30 mapping keys = 37 orphans"	65 disk vs 52 mapping entries = 13 orphans	verifier WP8
V2-P6	"verdict-ledger 100% force_completion"	79/107 rows (74%) force_completion; 28 standalone/done; PROBE_PASS=0 (gate-gaming concern stands)	verifier CP8
V2-P7	"claude-builder queue 2,945 failed / 1 completed"	TWO subsystems: queue-table has 2,944 rows (verifier WP3); durable-runner.db has 295/1/1 completed/failed/pending (Proveo C-04). MC-C4-2 NEEDS RE-PROBE before dispatch.	Proveo C-04 + verifier WP3
V2-P8	"TLDR daemon writes to ~/system/data/insights/ which does not exist"	Daemon writes to ~/system/logs/tldr-insights/ which EXISTS with files from 2026-04-24. MC-C1-4 scope needs re-audit.	Proveo C-11
V2-P9	"manifest-index.md last 2026-02-26"	mtime 2026-04-06 (Feb 26 is content audit date inside file); 43 days stale	verifier KP3
V2-P10	"HiveMind 21,741 rows"	21,930 live (audit-snapshot drift)	verifier KP5
V2-P11	"True 7d = \$365,104"	\$366,236 (Proveo C-10, ±0.3% rounding)	Proveo C-10
V2-P12	"MC backlog blocked = 2,239"	2,241 (Proveo C-02, +2 drift)	Proveo C-02

### Re-probe required (BLOCKERS for build dispatch):

- MC-C4-2 (claude-builder drain decision) — Team C must specify exact DB path + table before scope freeze
- MC-C1-4 (TLDR daemon fix) — re-audit actual writer path vs `~/system/logs/tldr-insights/`
- WP6 "2,400 zombie MCs" — verifier blocked by bash-danger-gate; needs read-only sqlite policy fix or alternate probe

**Verdict on v2.0 after patches:** Strategic narrative + 4-week roadmap + 9 CEO decisions HOLD. Six precision errors corrected in this section. v2.0 is publication-ready with footnoted re-probes on MC-C4-2 + MC-C1-4.

---

Revision #2

Created 2026-05-19 15:55:25 UTC by John

Updated 2026-06-21 20:03:50 UTC by John