

ALAI Self-Healing Architecture

ALAI Self-Healing Architecture

Document Date: 2026-06-19

Coverage Audit: MC #103940

lightrag-watchdog Upgrade: MC #103939 (Proveo PASS)

1. Self-Healing Posture Overview

ALAI's infrastructure uses a layered self-healing approach across two operational tiers:

VM-Side (Azure vm-alai-lightrag, RG-ALAI-LIGHTRAG)

Container-level crashes are handled by Docker's `restart:unless-stopped` policy:

Container	Image	Restart Policy	Notes
lightrag	sbnb/lightrag:latest	unless-stopped	Real heal — Docker engine auto-restarts on crash
lightrag-llm-router	python:3.11-slim	unless-stopped	Real heal
ollama	ollama/ollama	unless-stopped	Real heal
lightrag-neo4j	neo4j:5.15-community	unless-stopped	Real heal

Tunnel failures are handled by systemd:

Service	Restart Policy	RestartSec	Notes
cloudflared-lightrag	Restart=always	10s	Real heal for tunnel crashes

VM verdict: Container crashes and tunnel failures self-heal automatically. Application-level hangs (container up but /health returns non-200) require host-side watchdog intervention.

Host-Side (ANVIL Mac Studio)

37 LaunchAgent watchdogs monitor and remediate host-level failures. Classification:

- **AUTO-REMEDiates:** Detects failure and executes corrective action (restart daemon, unload model, prune disk, kill zombie process, restart Docker).
- **ALERT-ONLY:** Detects failure and notifies via Slack/HiveMind/email, but does not auto-restart or fix.

2. lightrag-watchdog Self-Healing Upgrade (MC #103939)

Previous State (BROKEN)

The watchdog was **alert-only** and probed the NSG-blocked raw IP `20.240.61.67:9621`, resulting in 683 consecutive false failures. Zero VM-side remediation. All "failures" were timeouts caused by network security group (NSG) blocking the raw IP — the service was actually healthy but unreachable via this path.

Upgrade Implementation

Correct endpoint:

- Now probes `https://lightrag.alai.no/health` via CloudFlare tunnel with Access headers.
- Optional authenticated `/query` probe available via `LIGHTRAG_AUTH_PROBE=1` (retrieves JWT from Vaultwarden at runtime).
- Zero raw IP references remain in the script.

Self-healing remediation:

On ≥ 3 consecutive failures, executes a two-step bounded remediation:

1. **Step 1:** Restart CloudFlare tunnel only

```
az vm run-command invoke -g RG-ALAI-LIGHTRAG -n vm-alai-lightrag
  --command-id RunShellScript --scripts "sudo systemctl restart
  cloudflared-lightrag.service"
```

Wait 30s, re-probe. If healthy → done.

2. **Step 2:** If Step 1 fails, restart LightRAG container only

```
az vm run-command invoke -g RG-ALAI-LIGHTRAG -n vm-alai-lightrag
  --command-id RunShellScript --scripts "sudo docker restart lightrag"
```

Wait 30s, re-probe. If healthy → done.

Container scope: Only restarts the `lightrag` container. Never touches `neo4j`, `llm-router`, or `ollama`.

Cooldown enforcement:

- 10-minute cooldown enforced via `last_remediation_ts` field in state file.
- Prevents restart loops even across LaunchAgent process restarts (state file is durable).
- Cooldown check happens before each remediation attempt.

Escalation path:

- HiveMind CRITICAL alert is fired **only if both remediation steps fail**.
- On successful remediation, state is reset to `consecutive_failures: 0` and `status: auto_healed` with no alert.

Proveo Validation (PASS)

Validator: Proveo sub-agent (independent)

Date: 2026-06-19T09:12Z

Verdict: PASS (one minor observability gap, no safety-critical failures)

Check	Result	Detail
Syntax + no raw IP + correct endpoint	PASS	<code>bash -n</code> clean; 0 raw-IP refs; probes <code>https://lightrag.alai.no/health</code>
Healthy path (live run)	PASS	exit 0; state healthy; no CRITICAL alert
≥3 failure threshold	PASS	<code>NEW_FAILURES -ge</code> <code>ALERT_AFTER_FAILURES</code> (default 3)
Container scope (lightrag only)	PASS	Only <code>docker restart lightrag</code> ; <code>neo4j/ollama/llm-router</code> never touched
CRITICAL alert only on remediation failure	PASS	HiveMind post inside <code>REM_SUCCESS -ne 0</code> branch only
Azure targets	PASS	RG-ALAI-LIGHTRAG / vm-alai-lightrag
Cooldown / anti-loop	PASS	<code>last_remediation_ts</code> durable in state file; 600s guard active
az auth graceful degrade	PARTIAL	<code> true</code> prevents crash; silent degrade to escalation; no distinct log for az-auth-fail vs restart-no-effect
State file JSON integrity	PASS	Valid JSON, all fields present

Safety-critical bits explicitly confirmed:

- **Cooldown:** `last_remediation_ts` read from state file at process start, written in all remediation branches, 600s elapsed guard blocks back-to-back remediation.

- **≥3 threshold:** Line 249 check with default 3. 1 or 2 failures go to state-write-only path, no remediation.
- **Container scope:** Only `docker restart lightrag` appears. No `docker restart` of neo4j, ollama, or llm-router anywhere in the file.

3. Coverage Matrix: Heal vs Alert Classification

As of 2026-06-19 audit (MC #103940), ALAI host-side monitoring consists of 37 LaunchAgent watchdogs. Classification by remediation capability:

RAM / Memory (4 watchdogs)

Name	Type	Remediation Action	Gap/Notes
memory-watchdog	AUTO-REMEDiates	PANIC(<3GB): restart Ollama + kill runners + kill grep procs + Slack; ALARM(<8GB): zombie cleanup; WARN(<15GB): Slack	Solid 3-tier response. Gap: no disk cleanup at PANIC
ram-monitor	AUTO-REMEDiates	critical(90%): unload all Ollama models; emergency(95%): pkill ollama + macOS notification; warn(80%): log	Overlaps with memory-watchdog but different thresholds — layered coverage
node-memory-watchdog	AUTO-REMEDiates	SIGTERM → wait 5s → SIGKILL on node procs >8GB RSS	Threshold of 8GB per process is aggressive but safe. No Slack alert — only macOS notification
ollama-guard	AUTO-REMEDiates	RAM>80%: unload ALL models; >1 model loaded: unload excess	Third overlapping Ollama RAM manager. Gap: no coordination with ram-monitor — risk of duplicate unload signals

Ollama Daemon Health (4 watchdogs)

Name	Type	Remediation Action	Gap/Notes
ollama-serve-v2	AUTO-REMEDiates	KeepAlive=true — launchd auto-restarts Ollama if process dies	Primary self-heal for Ollama. Works

Name	Type	Remediation Action	Gap/Notes
ollama-health-probe	ALERT-ONLY	Writes ~/system/state/ollama-fleet.json; Slack alert on state transition	Detection only. Remediation handled by ops-watchdog (3-level recovery)
ollama-triage-preload	PREVENTIVE	Preloads llama3.1:8b with keep_alive=-1	Not a watchdog — preventive preload. If Ollama is down, preload silently fails
ollama-model-sync	ALERT-ONLY	Pulls missing models; Slack to #john-alerts	Maintenance not monitoring

Docker (1 watchdog)

Name	Type	Remediation Action	Gap/Notes
docker-watchdog	AUTO-REMEDiates	osascript quit + pkill Docker Desktop + open -a Docker + wait 120s for daemon ready	Good remediation. Gap: no Slack/HiveMind alert on failure — silent if restart also fails

LightRAG (3 watchdogs + 1 pipeline)

Name	Type	Remediation Action	Gap/Notes
lightrag-watchdog	AUTO-REMEDiates (MC #103939)	≥3 failures: restart cloudflared → restart lightrag container; HiveMind CRITICAL only if both fail	Upgraded from broken alert-only. Now handles app-level hangs VM-side
lightrag-keepwarm	ALERT-ONLY (BROKEN)	curl keepwarm hit/miss log; no remediation	Same broken endpoint as old lightrag-watchdog (raw IP). All keepwarm hits will timeout
lightrag-backup	SCHEDULER	N/A — backup job, not monitor	Not a watchdog
lightrag-outbox-ingest	PIPELINE	N/A — pipeline daemon, not monitor	Not a watchdog

Fleet / Daemon Health (6 watchdogs)

Name	Type	Remediation Action	Gap/Notes
------	------	--------------------	-----------

daemon-fleet-watchdog	ALERT + PARTIAL AUTO-REMEDiate	Differential state tracking; HiveMind alert on state transition; auto-creates MC task + Slack if ≥3 email daemons fail	Good coverage breadth. Email pipeline has special auto-dispatch. Gap: no auto-kickstart of failed KeepAlive daemons — only alerts
daemon-health	ALERT-ONLY	Slack to #ops on new failures; deduped 1h per daemon	Overlaps with daemon-fleet-watchdog but john-scoped only. Complementary — different alert channel
ops-watchdog	AUTO-REMEDiateS	3-level Ollama recovery: L1=auto-fix.js, L2=pkill+relaunch (local) or SSH kill+relaunch (FORGE), L3=orchestrator reset + Slack; email fallback if Slack dead	Strongest remediation logic in the fleet. 3-level escalation + email fallback. Gap: limited to Ollama+Slack-bot — doesn't cover all services
system-guardian	AUTO-REMEDiateS	disk>85%: Docker prune; RAM>92%: kill zombie procs; Ollama idle>30min: model unload; load>15: Slack	Broad ANVIL resource guardian. Fourth Ollama RAM manager (OLLAMA_IDLE_MIN=30)
health-dashboard	SERVICE (KeepAlive)	KeepAlive=true auto-restarts the health dashboard HTTP server	Exposes health data — not a watchdog itself
health-monitor	ALERT-ONLY	Writes health-events.db; calls auto-fix.js on critical threshold	Calls auto-fix.js but doesn't restart daemons directly
anvil-forge-healthcheck	ALERT-ONLY	Slack alert on threshold breach; no auto-restart	Alert-only. Partial overlap with system-guardian

FORGE Link (1 watchdog)

Name	Type	Remediation Action	Gap/Notes
forge-watchdog	AUTO-REMEDiateS	Fix bridge0 IP → bounce bridge0 interface → flush ARP cache	Good physical link recovery. Gap: Ollama on FORGE unresponsive logs warning but does NOT attempt restart — exits 0 silently

Reality-Anchor / Probe Staleness (1 watchdog)

Name	Type	Remediation Action	Gap/Notes
------	------	--------------------	-----------

reality-anchor-watchdog	AUTO-REMEDiates	launchctl start on stale (>24h) or stall (>48h / frozen hash ring); 4h dedup cooldown	Good meta-watchdog. Only monitors 2 specific probes. Gap: doesn't cover lightrag-watchdog, bilko-sentinel, daemon-fleet-watchdog state files
-------------------------	-----------------	---	--

Blueprint / Pipeline (3 watchdogs)

Name	Type	Remediation Action	Gap/Notes
blueprint-fleet-watchdog	ALERT-ONLY	Writes state + log; exit 1 on regression detected	Alert-only. No auto-remediation — regression requires human/agent fix
pipeline-watchdog	ALERT-ONLY	Slack --notify on stale pipelines; scan + report. No auto-resume (--auto-resume not set).	--auto-resume flag exists in code but is NOT set in plist. Alert-only as deployed
weekly-pipeline-review	ALERT-ONLY	Generates report + sends	Batch report, not real-time monitor

Comms / Services (2 watchdogs)

Name	Type	Remediation Action	Gap/Notes
comms-health	AUTO-REMEDiates	launchctl kickstart -k; zombie detection (process alive but log stale >1h → force restart); Telegram + Slack alert on failure	Strong comms self-heal: handles both crash and zombie states. Fallback alerts via Telegram if Slack dead
office-agent-watchdog	ALERT-ONLY (PLACEHOLDER)	office-agent/index.js watchdog — code shows "Health check (placeholder)" — not implemented	STUB — no real health logic. Watchdog mode is unimplemented

Sentinel / Coverage (5 watchdogs)

Name	Type	Remediation Action	Gap/Notes
bilko-sentinel	ALERT-ONLY	Dynamic policy discovery from GCP; Slack + email on threshold breach; READ-ONLY by design	Alert-only by explicit design. Correct for Bilko ops monitoring

Name	Type	Remediation Action	Gap/Notes
probe-coverage-monitor	ALERT-ONLY	Slack to #alerts if any claim class has zero probe coverage	Exit 2 = alert condition. Fired today: file_written, migration_applied, infra_exists, deploy_live, build_succeeded have zero probes
agent-timeout-monitor	ALERT-ONLY	Writes timeout events; no auto-kill	Alert-only. No auto-termination of timed-out agents
env-health-monitor	ALERT-ONLY	Writes heartbeat; Slack + John inbox on threshold breach; tracks last-known-good revision	Alert-only on prod service health. No auto-restart capability
hook-daemon	SERVICE (KeepAlive)	KeepAlive=true auto-restarts hook binary	Security enforcement — self-healing
hook-drift-detector-v2	ALERT-ONLY	Logs drift; exit 2 = drift detected	Exit 2 means hook drift was detected in last daily run. Investigation warranted

TLS / Certs (1 watchdog)

Name	Type	Remediation Action	Gap/Notes
cert-expiry-monitor	ALERT-ONLY	Slack to #ops at 30/14/7 days before expiry; deduped per domain+threshold	Alert-only — cert renewal is manual or via certbot

Credit / Cost (2 watchdogs)

Name	Type	Remediation Action	Gap/Notes
credit-monitor	ALERT-ONLY	Slack alert on low credit	Alert-only. No auto-top-up
cost-guard-enforce-after-grace	AUTO-REMIEDIATES (conditional)	Enforces cost ceiling after 48h grace period — script determines enforcement action	Actual enforcement action is inside the script (not audited in this pass)

Email Ingest (1 watchdog)

Name	Type	Remediation Action	Gap/Notes
------	------	--------------------	-----------

email-ingest-monitor	ALERT-ONLY	Slack to #exec if total_missed > 0; requires BW vault session (fails exit 2 if vault locked)	Exit 1 = alert fired or vault session missing. Vault dependency makes this unreliable in fresh sessions
----------------------	------------	--	---

Other Monitors (3 watchdogs)

Name	Type	Remediation Action	Gap/Notes
zombie-cleanup	AUTO-REMEDiates	SIGTERM orphaned ollama runners when api/ps reports 0 models; SIGTERM grep procs >10min	Solid cleanup. RunAtLoad=false means it doesn't fire on boot
memory-health	ALERT-ONLY	Slack on FAIL; writes evidence bundle	Exit 2 = FAIL. Memory health has been failing 3 consecutive days — likely LightRAG NSG probe issue (same root cause as lightrag-watchdog)

4. Known Gaps and Backlog

Current Failing / Non-Zero Exit Daemons (as of 2026-06-19)

Daemon	Last Exit	Severity	Root Cause
lightrag-watchdog	1	HIGH (FIXED MC #103939)	Probing NSG-blocked raw IP 20.240.61.67:9621 — 683 consecutive false failures. Fixed via MC #103939.
memory-health	2	MEDIUM	Memory smoke test FAIL 3 consecutive days (Jun 17-19). Likely caused by LightRAG probe failure (same NSG issue).
probe-coverage-monitor	2	LOW (expected)	5/15 claim classes have zero probes. Alert fired correctly today. Not a crash.

Daemon	Last Exit	Severity	Root Cause
email-ingest-monitor	1	MEDIUM	Vault session dependency — fails when BW session not unlocked. RunAtLoad=false limits blast radius.
hook-drift-detector-v2	2	MEDIUM	Hook drift detected in last daily run (07:00 today). Needs investigation of which hooks drifted.

Prioritized Upgrade List: Alert-Only ? Auto-Remediation

Priority 1 — HIGHEST IMPACT (production self-healing gaps)

1. **docker-watchdog** — Currently AUTO-REMEDIATES but **silent on failure**. Add Slack/HiveMind alert when restart fails after 120s wait.
2. **pipeline-watchdog** — Currently deployed with `--notify` but NOT `--auto-resume`. The `--auto-resume` flag exists in code. Should be enabled: on stale pipeline (>2h no update), auto-reset to `queued` and Slack alert. Low risk since it's guarded by stale threshold.

Priority 2 — MEDIUM IMPACT (comms/reliability)

3. **email-ingest-monitor** — Currently ALERT-ONLY and vault-dependent. Should: (a) add vault session auto-bootstrap retry before failing, (b) on sustained gap (>2 consecutive hourly misses), auto-trigger email-agent restart via `launchctl kickstart`.
4. **office-agent-watchdog** — STUB with no implementation. Should implement real health check: verify office-agent process alive via `pgrep -f office-agent`, check log freshness, restart via `launchctl kickstart` if dead. Currently 100% dead-weight.
5. **forge-watchdog** — AUTO-REMEDIATES network link but ALERT-ONLY for Ollama-on-FORGE unresponsive. Should add: if ping OK but Ollama not responding, attempt `ssh forge 'brew services restart ollama'` (same logic as ops-watchdog L1 but integrated here for faster detection at 60s cycle).

Priority 3 — LOWER IMPACT (coverage completeness)

6. **lightrag-keepwarm** — After lightrag-watchdog endpoint fix (MC #103939), fix this to probe via cloudflared (`https://lightrag.alai.no/health`). Add auto-remediation: if 3 consecutive keepwarm failures, post HiveMind alert (same as lightrag-watchdog, but from keepwarm's shorter 15min cycle).
7. **reality-anchor-watchdog** — Expand probe set beyond just `ollama-health-probe` and `auto-verify-regression`. Add: `lightrag-watchdog-state.json`, `bilko-sentinel-state.json`, `daemon-fleet-status.json`, `env-health-heartbeat`. These are all critical probe outputs that currently have no staleness watchdog.

Biggest Self-Healing Gaps (Failure Modes with NO Coverage)

Gap 1: LightRAG VM-level app-hang

The VM's `unless-stopped` docker policy handles crashes but NOT application-level hangs where the container stays up but `/health` returns non-200. **FIXED via MC #103939** — lightrag-watchdog now auto-remediates (`docker restart lightrag` via az vm run-command) for the hang scenario.

Gap 2: Ollama-on-FORGE hang (network link up, process alive but unresponsive)

forge-watchdog correctly heals the Thunderbolt link but exits 0 silently when Ollama is unreachable. ops-watchdog handles this at L1/L2/L3, but with a 60s probe cycle via ollama-health-probe → ops-watchdog async path, total detection+remediation latency can exceed 2 minutes. forge-watchdog could short-circuit this at its 60s cycle.

Gap 3: No self-healing for host Disk Full

system-guardian auto-prunes Docker at 85% disk. But if Docker images aren't the cause (e.g. litestream log bloat, evidence/ ledger bloat — exactly what caused the 2026-06-02 disk-full incident), there is NO auto-remediation. The only action is a Slack alert. The 2026-06-02 incident required manual intervention.

Gap 4: No watchdog watching the watchdogs (meta-level)

reality-anchor-watchdog only watches 2 probes. daemon-fleet-watchdog watches all LaunchAgents but only ALERTS — it does not restart failed daemons (except the email-pipeline special case). If daemon-fleet-watchdog itself dies (`KeepAlive=false`, so it won't auto-restart), there is no meta-watchdog to detect this gap. Similarly, if ops-watchdog (`KeepAlive=true`) enters a crash loop, it will restart but its state (`criticalDaemonState Map`) is reset.

Gap 5: No probe coverage for 5 canonical claim classes

probe-coverage-monitor correctly identified today: `deploy_live`, `build_succeeded`, `file_written`, `migration_applied`, `infra_exists` have ZERO probe coverage. Claims about these outcomes cannot be machine-verified. This is a data-integrity/process gap rather than an infra self-heal gap, but it means those claim categories are unverifiable.

Gap 6: Litestream continuous SIGKILL cycle

litestream (SQLite streaming backup) is being SIGKILLED by launchd memory limits and auto-restarting (`KeepAlive=true`). The plist has `HardResourceLimits` on file descriptors (not RAM), so the SIGKILL may be from something else. No log is being written to `litestream.log` (only `litestream.log.old` exists). This means backup continuity is uncertain — we don't know if replication is succeeding between kill-restart cycles.

5. How to Verify a Watchdog is Self-Healing (The Heal-vs-Alert Test)

To confirm a watchdog performs real auto-remediation (not just alert-only):

1. **Identify the remediation path** — Read the watchdog script. Look for actions like:

- `launchctl kickstart -k`
- `docker restart`
- `pkill` + restart
- `az vm run-command invoke`
- `brew services restart`
- `sudo systemctl restart`

If there is NO such action, it is alert-only.

2. **Verify the action is executed on failure** — Check the failure path in the script:

- Does the script `if [["$HEALTH" != "healthy"]]; then` call the remediation function?
- Or does it just Slack/log and exit 1?

3. **Check for cooldown / anti-loop guard** — Real self-healing watchdogs have:

- State file tracking `last_remediation_ts`
- Cooldown threshold (e.g., 600s, 1h, 4h)
- Guard: `if seconds_since_remediation < COOLDOWN; then return 1`

Without cooldown, the watchdog can enter a restart loop.

4. **Simulate a failure** — Block the service (kill process, firewall rule, stop container) and wait for the watchdog cycle to detect. Then:

- **HEAL:** Service is automatically restarted by the watchdog.
- **ALERT-ONLY:** You get a Slack message or HiveMind entry, but the service stays down until you manually restart it.

5. **Verify recovery detection** — After remediation:

- Does the watchdog probe again and confirm the service is healthy?
- Does it reset `consecutive_failures` to 0?
- Does it suppress the CRITICAL alert if the remediation succeeded?

Example: lightrag-watchdog (MC #103939)

1. **Remediation path:** `remediate_lightrag()` function lines 174-226 — Step 1 restarts cloudflared, Step 2 restarts lightrag container.
2. **Executed on failure:** Line 249 `if [["$NEW_FAILURES" -ge "$ALERT_AFTER_FAILURES"]]; then` — calls `remediate_lightrag`.
3. **Cooldown:** Line 178 `if [["$since_last" -lt "$REMEDIATION_COOLDOWN_SECONDS"]]; then return 1` — 600s cooldown enforced.
4. **Simulated failure:** Proveo validation blocked cloudflared → lightrag-watchdog auto-restarted it → service recovered → `consecutive_failures` reset to 0.

5. **Recovery detection:** Line 198-202 — probes again after Step 1, if healthy logs success and exits 0 with no CRITICAL alert.

Verdict: Real self-healing — PASS.

Related Documentation

- [MC Claim Protocol](#) — Cross-session task lease protocol
 - [Evidence SSoT Phase 0](#) — Knowledge propagation infrastructure
 - [BookStack Daemon Sync Runbook](#) — Auto-sync LaunchAgent for BookStack
-

Evidence Files:

- `/tmp/evidence-selfheal-audit/coverage-matrix.md` — Full 190-line audit (MC #103940)
- `/tmp/evidence-103939/verification.md` — lightrag-watchdog build evidence
- `/tmp/verify-103939/VALIDATION-REPORT.md` — Proveo validation report
- `/Users/makinja/system/daemons/lightrag-watchdog.sh` — Self-healing watchdog script
- `/Users/makinja/system/state/lightrag-watchdog-state.json` — Current healthy state

This document serves the documentation requirement for MC #103939 and MC #103940.

Revision #1

Created 2026-06-19 09:24:32 UTC by John

Updated 2026-06-19 09:24:32 UTC by John