

REPRODUCIBILITY DOCUMENTATION – EVERY CITED NUMBER, INDEPENDENTLY VERIFIABLE

Lemur Audit – Raw Metadata

Every number that can be cited about the audit, with how it was recorded and how to verify it. For press-kit fact-checking, journalist due diligence, and any defensible-claim build.

1. The run

An audit of the Netflix/Lemur repository at pinned commit [f478458](#) (Lemur 1.9.0), produced by Umma in a single autonomous run. The final artifact is a 52-sheet XLSX workbook (1.29 MB) – status: ready, version 1.

2. Time

| THING | VALUE |
|-------------------|---|
| Run started | 2026-05-29 18:02:25 UTC |
| Run completed | 2026-05-29 21:22:32 UTC |
| Total duration | 3 hours 20 minutes 7 seconds |
| Iteration history | 28 conversations between 2026-05-26 and 2026-05-29 – the prompt was converged across 4 days |

3. Decomposition

| THING | VALUE |
|------------------------------|--|
| Root thoughts | 1 |
| Child thoughts | 36 |
| Completed child thoughts | 35 (97%) |
| Distinct decomposition roles | 9 |
| Tree depth | 2 (root + one level of children – flat by design for a parallel fan-out audit) |

3.1 Children by role

| ROLE | COUNT |
|------------|-------|
| implement | 8 |
| synthesize | 7 |
| research | 6 |
| evaluate | 5 |
| decide | 3 |
| verify | 3 |
| review | 2 |
| code | 1 |
| design | 1 |

4. Tool invocations

| THING | VALUE |
|---------------------|-----------|
| Total invocations | 907 |
| Successful | 635 (70%) |
| Failed / retried | 272 (30%) |
| Distinct tool names | 25 |
| Average duration | 505 ms |

4.1 Top tools by call count

| TOOL | CALLS |
|---|-------|
| <code>read_github_file</code> | 327 |
| <code>execute_code</code> | 220 |
| <code>repl_execute</code> | 192 |
| <code>get_code_file</code> | 32 |
| <code>find_existing_code_project</code> | 26 |
| <code>search_web</code> | 17 |
| <code>create_code_project</code> | 17 |
| <code>create_code_file</code> | 13 |
| <code>start_code_change</code> | 12 |
| <code>list_code_projects</code> | 11 |

(Full 25-tool list available on request.)

5. Cognitive layer outputs

| LAYER | COUNT |
|-------------------------|-------|
| L1 interpretations | 69 |
| L2 syntheses | 7 |
| Provenance links | 71 |
| Final artifact sections | 52 |

5.1 Interpretations by module

| MODULE | INTERPRETATIONS | AVG CONFIDENCE |
|------------|-----------------|----------------|
| umma | 33 | 1.00 |
| implement | 8 | 0.75 |
| synthesize | 7 | 1.00 |
| research | 6 | 1.00 |
| evaluate | 5 | 1.00 |
| decide | 3 | 0.67 |
| verify | 3 | 1.00 |
| review | 2 | 1.00 |
| design | 1 | 1.00 |
| code | 1 | 1.00 |

Why this matters for credibility: the lower confidence on `implement` (0.75) and `decide` (0.67) is what calibrated, honest reporting looks like. Confident assertion on every module would be a red flag for LLM-output overconfidence.

5.2 Headline synthesis

| FIELD | VALUE |
|--------------------------|-------|
| Confidence overall | 0.917 |
| Convergences | 36 |
| Divergences | 41 |
| Unresolved disagreements | 0 |

6. Artifact composition

| THING | VALUE |
|--|--|
| Final XLSX sections (sheets) | 52 |
| Total structured content | 1,288,348 bytes (1.29 MB) |
| Master findings rows | 238 |
| Master findings columns | 36 |
| Canonical findings (post-reconciliation) | 236 |
| Severity tier distribution | 52 T1 / 123 T2 / 46 T3 / 15 T4 |
| Named compounding chains | 8 (6 surfaced in headline CC sheet) |
| Compounding chain tier distribution | 7 T1 / 1 T2 |
| API endpoints mapped | 65 |
| Trust boundaries identified | 14 |
| Threat actors modeled | 10 |
| STRIDE matrix rows | 26 |
| Assets catalogued | 17 |
| Hardening matrix rows | 33 (operator-lever) + 5 (structural-residual) = 38 |
| Embargo redactions applied | 22 |
| Invariant checks passed pre-handoff | 24 |
| Files inventoried in codebase | 575 |
| Plugins registered (setup.py entry_points) | 33 |
| Plugins existing but not registered | 2 |
| Lemur native GHSAs cleared at pin | 4 |
| Public-research correlations | 10 |
| Python dependency findings | 29 |
| Frontend security findings | 4 |
| Unfixed dependency CVEs identified | 3 (urllib3, AngularJS, Bootstrap) |
| Per-endpoint rollup rows | 19 |
| Reconciliations applied | 5 final + 5 reconciliation logs |
| OCR ambiguities recorded | 3 |

7. Verifiability

Every figure above is recorded in the run's structured provenance — the decomposition tree, the per-thought tool log, the synthesis record (with its convergences, divergences, and unresolved-disagreement counts), and the final artifact's 52 sheets — and can be independently re-derived from the artifact and its provenance. Query-level access for fact-checking is available to press and partners on request.

8. What this audit did and did not do

For the press kit, it's worth being precise about both sides.

8.1 What the audit did

- Read **327 source files** from the Netflix/Lemur repository at the pinned commit ([f478458](#) , Lemur 1.9.0)
- Ran **220 sandboxed code-analysis tasks** — Python scripts that parsed Lemur's source (AST walks), traced call graphs across files, computed reachability between findings, ran custom analyzers (including a Semgrep-based analyzer), built the file inventory and endpoint map, and assembled the final XLSX
- Ran **192 interactive REPL sessions** for exploratory queries against the codebase
- Performed **17 web searches** to correlate against public security research, CVE databases, and GHSA records
- Browsed **4 web pages** and scraped documentation for cross-referencing
- Created **17 intermediate analysis projects + 13 files** for working state
- Built a full **STRIDE threat matrix (26 rows)** via systematic threat enumeration applied to each Lemur subsystem
- Constructed **6 named compounding attack chains** by composing findings across specialists — static cross-file flow analysis identifying chains where one finding enables another
- Built a **33-row deployment hardening matrix** by analyzing the public docker-compose configuration
- Identified the dependency findings by parsing requirement files, looking up affected version ranges in GHSA/CVE databases, and filtering by reachability into the call graph

These are real cognitive operations executed against the codebase. The **412 total code executions** (220 + 192) were Umma's own analysis tooling running in sandboxes — not Lemur being run.

8.2 What the audit did NOT do

The distinction is between running analysis code against the source (§8.1) and running the Lemur application itself:

1. **Did not run Lemur as a live service or send exploit payloads.** The compounding chains and per-finding severity scores are static-analysis reasoning artifacts — they identify reachability and attack composition by reading the code, not by sending HTTP requests to a running Lemur.

Confirming any specific chain end-to-end would require dynamic testing, which was not performed.

2. **Did not perform dynamic / runtime / fuzz testing.** Race conditions, real RBAC enforcement under concurrent load, runtime-specific deserialization gadgets, plugin-interaction bugs that manifest only at runtime – out of scope. Static analysis catches structural weaknesses; dynamic testing catches behavioral ones. The audit did the first, not the second.
3. **Did not validate patches via runtime regression testing.** When a CVE is “fixed at pin,” the audit confirms the fix is present in the source at the pinned commit but does not run the patched code under attack to confirm completeness.
4. **Did not perform OSINT on Netflix’s production deployment.** Findings about docker-compose defaults are based on the public repository’s deployment artifacts, not Netflix’s actual infrastructure, which may have additional hardening.
5. **Did not perform supply-chain forensics.** The dependency findings identify vulnerable versions (urllib3, AngularJS, Bootstrap) but do not investigate whether the published dependencies were tampered with. Supply-chain integrity verification is its own discipline and out of scope.

8.3 What this means for the press claim

The audit is what it claims to be: integrated autonomous production of a senior-consultancy-grade security audit in 3 hours 20 minutes, derived from sustained static analysis backed by real code execution against the source. It is *not* a penetration test of Netflix’s production infrastructure; nobody should read it as that. The recorded trail makes the distinction precise: 220 sandboxed analysis executions, zero requests sent to a live Lemur instance. Both honesty notes – what the audit did and what it didn’t – are deliberate scope choices, not oversights. The combination is part of what makes the rest credible.