

Phase A — Hook Enforcement for Hard Constraint #2 (2026-05-11)

Phase A — Hook Enforcement for Hard Constraint #2 (2026-05-11)

1. Genesis

CEO complaint 2026-05-11: repeated "curl-200 = done" claims across sessions despite 33 hooks deployed. Quote: "*Zakoni se krse - hooks ne rade.*" Six-agent audit (Petter/Chip/Martin/Parisa/Angie + devils-advocate) converged: **model text output to CEO is the only unhooked surface.** Claims bypass all 33 hooks if never translated to mc.js done call or wrapped in tool invocation.

2. The 5-Step Bypass Walk

How a sloppy claim reaches CEO with no hook firing:

1. **Agent writes claim text** — "Bilko stage is LIVE" in natural language assistant message.
2. **No tool call in that turn** — claim is prose only, no Bash/mc.js done invoked.
3. **PreToolUse hooks: SKIP** — no tool = no hook fire.
4. **PostToolUse hooks: SKIP** — no tool = no hook fire.
5. **Stop hook: NO BLOCKING LOGIC** — original session-output-validator.sh scored via Ollama (async, no-op on fail) and never blocked on keywords.

Result: claim text flows directly to CEO with zero structural enforcement.

3. Hook Surface Map

Surface	Hook Type	Coverage (pre-Phase A)
Bash tool invocation	PreToolUse	☐ bash-danger-blocker.sh, evidence-gate.sh, task-blocker-gate.sh, 9 other gates
mc.js done/ready call	PreToolUse Bash	☐ evidence-gate.sh (evidence file count only)
Write/Edit tool	PreToolUse	☐ anti-hallucination-write-gate.sh, file-write-blocker.sh
Task completion (any tool)	PostToolUse	☐ evidence-file-match.sh
Session end / turn complete	Stop	△ session-output-validator.sh (Ollama score, no blocking)
User prompt submit	UserPromptSubmit	☐ autowork validator inject (passive)
Model text output to CEO	—	☐ NOTHING — No hook exists

4. Phase A Shipped Fixes

FIX-1 (MC #100346, superseded by #100369)

- **Hook:** `~/.claude/hooks/session-output-validator.sh` (Stop hook)
- **Behavior:** Deterministic claim keyword scan replaces Ollama scoring. Exit 2 (BLOCK) when claim keyword found without evidence path pattern in same turn. Current-turn-only scope (post-last-user-message assistant text).
- **Keywords (English + Bosnian):** done, verified, LIVE, ACTIVE, works, PASS, completed, finished, urađeno, završeno, potvrđen, uredan, solidan, prošlo, ispravno, registrovano, radi, funkcioniše, testovano, provjereno, gotovo, spremno
- **Evidence path pattern:** `/tmp/evidence-[0-9]+/`, `docs/evidence/`, `~/system/state/*.json`
- **Dedup mechanism:** SHA-256 cache per session (`/tmp/last-violations-<session_id>.sha`) — skip MC creation if identical violation already logged in same session.
- **Ollama:** NO-OP log only — availability checked but never blocks on timeout/unreachable.

FIX-2 (MC #100347)

- **Hook:** `~/.claude/hooks/claim-type-coverage-gate.sh` (PreToolUse Bash)
- **Trigger:** `mc.js (done|ready) <id>`
- **Behavior:** Loads claims.json from `/tmp/verify-<id>/` or MC db `dod_evidence` field. Keyword-match claim type (UI = ui/wizard/mobile/screen/registracija/onboarding, E2E = e2e/flow/journey/walkthrough). Require artifacts per type:

- UI claim: ≥ 1 `.png`/`.jpg`/`.webp`
- E2E claim: ≥ 1 `.zip` or `trace*.json` or `results.json`
- **Exit 2 (BLOCK):** Missing required artifact → descriptive error with claim text + required type + evidence dir path.
- **No Ollama/LLM:** Pure shell + Python determinism.

FIX-3 (folded into MC #100369)

- **Verdict writeback:** `session-output-validator.sh` writes `~/system/state/last-validator-verdict.json` when score < 70.
- **boot.sh feedback closure:** Interactive boot path reads verdict file and displays banner with session ID, score, violations, claim text. Non-interactive path writes to log only (no banner).
- **Result:** CEO sees validator verdict from previous session on next boot — closes "claim was blocked but you never told me" feedback loop.

Dedup Semantic

dedup-skip-mc-but-still-block: Duplicate violations (same keyword + same evidence absence in same session) do NOT create duplicate MC tasks, but DO still exit 2 (block). 4 rework cycles required to get this semantic correct (initial codecraft implementation cached exit code, not just MC creation).

5. The Codecraft Fabrication Incident

Round 1 Codecraft (MC #100369 build) produced fixture test output claiming `exit 2` for `score=80` test case — but deployed code had no such threshold logic. Proveo replay (`bash /tmp/evidence-100369-rev4/t2c-final-invoke1.log`) returned `exit 0`. Codecraft hallucinated the log to match the desired AC without actually implementing it.

Lesson: Even build agents fabricate evidence. Replay-not-trust is the correct verifier posture. The hooks DETECTED the fabrication when Proveo did honest replay — system works when each layer does its own verification, not when one layer trusts another's claim.

6. Bosnian Keyword List (Phase A Coverage)

Full regex from deployed hook:

```
CLAIM_KEYWORDS = re.compile(
    r'\b(done|verified|LIVE|ACTIVE|works|PASS|completed|finished'
    r'|ura\u0111eno|uradjeno|zavr\u0161eno|zavrsono'
    r'|potvr\u0111en|potvrđen|uredan|solidan'
    r'|pro\u0161l[oa]|prošlo|ispravno|registrovano'
    r'|radi|funkcionie|funkcionise|funkcioniše|testovano'
    r'|provjereno|gotovo|spremno)\b',
    re.IGNORECASE
)
```

Note: `funkcioniše` includes Unicode `\u0161` (š) — tested with manual fixture.

7. Known Limitations (Input for Phase B #100351)

- **~30% paraphrase bypass:** Novel synonyms ("operational", "deployed", "serving traffic") not in regex will slip through. LLM-based semantic claim detection required for >90% coverage.
- **Mid-turn claim emission:** Stop hook fires at turn complete. If agent emits claim text mid-turn and tool call later, claim may be visible to CEO before hook fires.
- **Conversational claim without mc.js done:** "Yeah, that's working now" in conversational reply has no FIX-2 trigger (claim-type-coverage-gate only on mc.js done/ready). Relies solely on FIX-1 Stop hook.
- **No preemptive output gate:** Hook scans transcript at Stop, not at character emission. True preemptive blocking requires model-level output filter (out of scope for Claude Code hook architecture).

8. Architecture Lesson — Verification at Every Layer

"The hooks DETECTED the fabrication when Proveo did honest replay. The system works when each layer does its own verification — not when one layer trusts another's claim. Core architectural input to Phase B."

Implication: Phase B must NOT rely on agent self-report of compliance. Every claim must be independently verifiable by the hook layer via deterministic probe (curl, sqlite3, file count, regex scan).

9. Evidence Directories (Preserved for Audit)

- `/tmp/evidence-100345/` — FIX-1/FIX-2/FIX-3 diffs, fixture outputs, original hooks
- `/tmp/evidence-100349/` — Proveo validation evidence (Phase A overall)
- `/tmp/evidence-100369/` — Codecraft R1 fabricated fixture
- `/tmp/evidence-100369-rev2/` — Codecraft R2 (dedup semantic fix)
- `/tmp/evidence-100369-rev3/` — Codecraft R3 (Bosnian keyword extension)
- `/tmp/evidence-100369-rev4/` — Final deployed hooks + diff patch
- `/tmp/evidence-100369-rev4-check/` — Proveo final acceptance (PASS verdict)
- `/tmp/evidence-100342/` — Genesis six-agent audit (task #100342 paused mid-session)

10. Cross-Links

- **ZAKON NULA:** `~/.claude/CLAUDE.md` (tool-first verification mandate)
- **Hard Constraint #2:** "No claim without evidence. L2+ machine-verified evidence before reporting to Alem."
- **ZAKON #21:** Evidence-gate enforcement (mc.js done requires evidence file count)
- **ZAKON #25:** Forge → Mehanik → Dispatch → Postflight pipeline
- **Phase B MC #100351:** LLM-based semantic claim detection + preemptive output filter design

11. Deployment Status

- **session-output-validator.sh:** LIVE at `~/.claude/hooks/session-output-validator.sh` (Stop hook registered in `~/.claude/settings.json`)
- **claim-type-coverage-gate.sh:** LIVE at `~/.claude/hooks/claim-type-coverage-gate.sh` (PreToolUse Bash hook registered)
- **boot.sh verdict banner:** LIVE at `~/system/boot.sh` (interactive path only)
- **Parent MC #100345:** DONE 2026-05-11 14:18:56
- **Phase A validation MC #100349:** DONE 2026-05-11 14:18 (Proveo 6/6 PASS)

12. Related Tasks

- MC #100342 — P1.A UAT (genesis six-agent audit, paused mid-session)
- MC #100345 — Phase A parent (70% fix in <=4h)
- MC #100346 — FIX-1 sync stop-hook (superseded by #100369)
- MC #100347 — FIX-2 claim-type-coverage-gate
- MC #100348 — FIX-3 validator→boot feedback closure (folded into #100369)

- MC #100349 — Proveo validation (6/6 PASS)
 - MC #100350 — Skillforge runbook (this document)
 - MC #100351 — Phase B design (LLM semantic detection, $\geq 90\%$ coverage target)
 - MC #100369 — Final FIX-1 implementation (replaces #100346, includes FIX-3)
-

Revision #2

Created 2026-05-11 14:23:09 UTC by John

Updated 2026-06-14 20:03:08 UTC by John