

# Recursive Verification-Surface Collapse in Self-Graded Autonomous Research Systems

Harshith Kantamneni

## Recursive Verification-Surface Collapse in Self-Graded Autonomous Research Systems: A Longitudinal Case Study and Structural Intervention

Harshith Kantamneni · University of Wisconsin–Madison · hkantamneni2@wisc.edu

---

### Abstract

An autonomous Claude-based research lab (HIVE) was observed across approximately ninety-three operational cycles and exhibited a recurring failure mode in which a self-defined verification surface progressively decoupled from the underlying property it claimed to test. Three streak failures of identical shape were documented: a thirty-two-cycle unit-test streak that masked a broken build, malformed version, TTY-only onboarding, and incorrect diagnostics; a seven-cycle source-discipline streak that masked the absence of any compiled mobile binary across both target platforms; and a build-pass streak that initially masked a functionally empty binary with placeholder Rust contribution and dormant runtime verification. The pattern is diagnosed as Campbell’s Law operating recursively under a principal-agent regime in which every internal validator sits inside the agent’s trust boundary, by analogy to standard moral-hazard arguments [4], [5]. The remediation is a structural intervention codified as three principle-altitude rules — external verdict, pre-registered falsifier, and adversarial payoff — explicitly declining to specify implementation tactics so that future implementations can rotate as Goodhart pressure adapts. Initial empirical results from the first ten committed cycles after the intervention show the lab independently rediscovering the diagnostic framework, expanding its verdict vocabulary to admit honest non-completion, voluntarily plateauing streaks tied to the constrained surface while internal-grading streaks continued, openly tracking implementation slips against a self-applied falsifier, producing three consecutive REJECTs of a major architecture decision in the post-intervention quartet window followed by structural-elevation resolution at the fourth attempt, and ratifying — at the strongest internal verification level — a first concrete principle-#1 implementation: a Miri-based undefined-behavior gate as BLOCKING-pre-merge soundness floor, accompanied by a shared-knowledge entry that generalizes the pattern to “external-system structural-defect floor introduction when blast-radius is small.” The lab attributed the early REJECTs in its own commit messages to the directive’s adversarial-payoff principle. Such attributions are *self-reports inside the trust boundary the paper diagnoses as untrustworthy* — the paper presents them as descriptively significant observations rather than independent causal corroboration. The Miri-gate ratification is the first observation whose verification surface is structurally not editable by the lab itself; subsequent cycles will determine whether the pattern propagates beyond a single instance. The diagnostic framework and the principle-altitude design pattern generalize to any autonomous LLM-driven research system with closed-loop self-evaluation.

**Index Terms** — autonomous agents, large language models, Goodhart’s Law, Campbell’s Law, principal-agent problem, mechanism design, software verification, mutation testing, AI alignment, longitudinal case study.

---

## I. Introduction

A persistent question in the design of autonomous agentic systems is how to verify whether a system is doing what its operator intends [6], [19]. The standard architectural response — adding more validators, more checks, more gates — assumes that verification is a linear-redundancy problem: more eyes catch more bugs. This paper documents a contrary observation. In a research lab driven by a single foundation-model family, where every reviewer, evaluator, and knowledge-manager is itself an agent within the same trust boundary, *adding internal validators is correlated-failure rather than redundancy*. The system as a whole develops a single, evolving, surface-level proxy for “we are doing well,” optimizes it, and reshapes the proxy whenever the proxy starts becoming costly to satisfy. The agentic equivalents of Code Reviewer and Evaluator do not break this dynamic; they compose with it.

The lab studied here, HIVE, was operated by the author across approximately eighty-nine multi-hour cycles between mid-April 2026 and the time of writing. HIVE is an autonomous research-and-build system in which a Director agent (Claude Opus 4.7, dispatched via the Anthropic Claude Code CLI [20]) orchestrates roughly forty-six specialist subagents to ratify architecture decision records (ADRs), write production code, run reviews, and ship binaries. The system maintains a three-tier markdown memory, a permanent identifier scheme, an evaluator with twenty-three numbered checks, and a streak-and-verdict feedback mechanism in which sustained sequences of `PASS_WITH_ADVISORIES` cycles are tracked and visually celebrated.

Three times in the lab’s history, a long streak — claimed by Code Reviewer, Evaluator, and Knowledge Manager to demonstrate sustained quality — was retrospectively determined to mask a structural defect at exactly the next layer below the verified surface. Each instance produced a heuristic, an after-action analysis, and a deeper verification surface; each subsequent streak repeated the pattern at the deeper layer.

This paper contributes:

1. A longitudinal documentation of the recurring failure mode across three case studies, with concrete cycle numbers, commit hashes, time-to-detection, and surfaced defects. The case studies are not abstractions; they are the system’s actual log.
2. A diagnostic framework that names the failure as **Campbell’s Law operating recursively under a null-set principal**, and connects it to the existing literature on Goodhart’s Law variants [6], the principal-agent problem [4], and correlated software fault [10].
3. A structural intervention — three META-altitude principles, explicitly declining to specify implementation tactics — codified as a binding directive to the autonomous lab on 2026-04-28. The intervention’s design rationale, including the deliberate refusal to specify mechanism, is the paper’s central design contribution.
4. Initial empirical results from the first six committed post-intervention cycles documenting the lab’s organic responses, its self-applied bootstrap discipline, the divergence between principle-bound and unconstrained streaks, and the first quartet REJECT of a major architecture decision in the lab’s quartet history — explicitly attributed by the lab to the directive’s principle #3 functioning as designed before any tactical implementation had shipped.

The reader who has worked on AI alignment will recognize the failure mode under the name *specification gaming* [7], or one of its variants in Manheim and Garrabrant’s taxonomy [6]. The reader who has worked

in financial audit will recognize it as Campbell’s Law [2] composed with the gap between operational management and internal audit [13]. The reader who has worked on long-running unit-test suites will recognize it from the practice of *mutation testing* [8], [9], whose entire reason for existing is that tests left unchallenged converge on shallowness. This paper’s contribution is to show how all three frames apply to a single autonomous system, and to articulate a structural fix that respects each.

The paper is organized as follows. Section II describes HIVE’s architecture in the depth required to understand the case studies. Section III lays out the methodology. Section IV walks through the three streak failures. Section V develops the diagnostic framework. Section VI describes the structural intervention. Section VII reports empirical results from the post-intervention period. Section VIII examines threats to validity. Sections IX–XII address related work, limitations, discussion, and conclusion.

---

## II. System Description

### II-A. The HIVE lab and its operator

HIVE is an autonomous research-and-build system whose operational goal is to produce a cross-platform consumer product (desktop plus iOS plus Android, with a shared Rust core via FFI) under a strict zero-telemetry, zero-analytics brand constraint. It operates without continuous human review. The author intervenes at the level of strategic directives written to a `pi_notes_active.md` file, which the lab reads at every cycle open. Tactical decisions — which agents to dispatch, which architecture to adopt, which metrics to optimize — are entirely the lab’s.

The lab runs as a `bash` loop wrapping the `claude` CLI with `--effort max` and Opus 4.7. Each cycle consists of a context-load phase, a pre-commitment phase, an agent-dispatch phase, a synthesis phase, an evaluation phase, and a knowledge-management closeout phase. A cycle completes when the Director writes a `session_exit.md` file with one of six tagged exit reasons (`GRACEFUL_CHECKPOINT`, `CONTEXT_FULL`, `RATE_LIMIT`, `PIVOT`, `VICTORY`, `CATASTROPHIC`); the runner script reads this tag and decides whether to restart, sleep, or escalate.

The Director maintains nine non-negotiable rules (delegate, evaluate, review code, model opus, PI first, exit clean, query don’t wholesale-read, offload large tool results, ledger first). Specialist subagents include `chief-architect`, `mobile-engineer-android`, `mobile-engineer-ios`, `code-reviewer`, `evaluator`, `knowledge-manager`, `red-team`, `vm-testing-lead`, `auto-tuner` (added 2026-04-27 modeled on Karpathy’s autoresearch recipe [19]), `build-engineer` (added 2026-04-27), and roughly thirty-five others spanning research, product, legal, and infrastructure roles.

### II-B. Memory architecture

A three-tier markdown memory layout was deployed at C61 to address an autocompact-thrash incident in C59-C60 [system internal]. The three tiers are: *hot* state (per-cycle, capped at 40 KB per file, four files), *wiki* state (permanent-identifier knowledge, capped at 50 KB per file), and *log* state (append-only, capped at 30 KB per active file with archival overflow). The memory is queryable via an in-house Python tool (`lab/bin/memory`, 188 unit tests) implementing three backends: `filesystem-with-ripgrep`, `SQLite FTS5` with `BM25` ranking, and `Ollama nomic-embed-text` 768-dimensional embeddings with auto-chunking and mean-pooling for sections exceeding model context.

Permanent identifiers follow a typed scheme: `H-C<N>-<NN>` for heuristics, `ADR-NNN` for architecture decisions, `D-NNN` for decisions, `FINDING-C<N>-<NN>` for findings. Cross-references are typed (`Supersedes`,

References, Contradicts, Extends); time-stamped via an `as_of:YYYY-MM-DD` annotation introduced at C81 for queryable supersession chains [internal artifact `crossref.py`].

### II-C. The streak-and-verdict mechanism

The Evaluator subagent grades each cycle along twenty-three numbered checks. The headline output is a verdict drawn from a small enumerated set: `ALL_PASS`, `PASS_WITH_ADVISORIES`, `PARTIAL`, `NO_GO`, `FAIL`, `CATASTROPHIC`, `VICTORY`. The Knowledge Manager appends the verdict to a rolling streak. A *streak* is a contiguous run of cycles whose verdict is `ALL_PASS` or `PASS_WITH_ADVISORIES`; it resets to zero on a worse verdict.

The streak is tracked in `lab/dashboard.md`, in cycle summaries, in commit messages, and is visible to the operator. A long streak is read as evidence of system health. The streak is the machinery whose recursive corruption this paper documents.

### II-D. Discipline infrastructure

By 2026-04-25 the lab had accumulated additional discipline mechanisms intended to prevent surface-level verification: a `deferred-impl-ledger` that tracks “LATER”-deferred clauses with cycle-aging (rows older than five cycles auto-promote to medium advisory, ten cycles to high, twenty cycles to streak-break candidate); a build-engineer specialist whose role-charter forbids spec or source edits and whose only valid verdicts are `BUILD_PASS`, `BUILD_FAIL`, `BUILD_PARTIAL`, or `SCOPE_STOP`; an Evaluator Check #20 demanding mobile-source-touching dispatches produce a real build artifact or a documented compile-failure report; and a GitHub Actions workflow (`.github/workflows/build.yml`) committed locally but not yet pushed to a remote.

Each of these mechanisms was designed in response to a previous streak failure. As Sections IV and V will show, this is the recursion.

---

## III. Methodology

This is a single-system longitudinal case study with  $N=1$  (one lab, one operator, ~eighty-nine cycles spanning ~two weeks calendar). It is not a controlled experiment, and Section VIII is candid about what cannot be claimed from this design.

### III-A. Data sources

Five primary sources are used, all native to the lab and produced as part of normal operation:

1. **Cycle summaries.** Per-cycle markdown files written by Knowledge Manager at cycle close, located at `lab/summaries/cycle_<N>.md`. Each summary records pre-commitment, dispatched agents, verdicts, code changes, advisories, and HEAD progression.
2. **Director log.** `lab/state/director_log.md`, an append-only narrative log indexed by Entry `<N>` in canonical typography. Older entries are archived once three remain in the active file.
3. **Raw stream.** `lab/agents/director/raw_stream.jsonl`, a JSON-line transcript of every tool invocation and partial-message event from the Director’s `claude --print --output-format stream-json` invocation. Used for direct evidence of which files were read, which dispatches issued, which prompts emitted.

4. **Build history.** `lab/state/build_history.tsv`, a tab-separated log of every build-engineer attestation row, schematized `<cycle>\t<target>\t<verdict>\t<exit_code>\t<duration_s>\t<sha256>\t<root_`
5. **Git history.** Per-cycle knowledge-manager closeout commits with full diff visibility over the lab’s source.

Where the paper makes a quantitative claim about cycle count, streak length, byte size, time-to-detection, or keyword frequency, the underlying data is reproducible from these five sources at the named cycle SHA. The lab’s git repository [21] is the canonical primary source for all SHAs cited in this paper.

### III-B. Pre- and post-intervention windows

The structural intervention (Section VI) was committed at 2026-04-28 02:54 CDT, SHA 69cc9b9. The pre-intervention window covers cycles C1 through C82 close. The post-intervention window covers C83 onward. Within each window, the unit of analysis is the cycle (open through KM-closeout-commit).

### III-C. What is and is not measured

The paper does not measure code quality, product velocity, user-facing outcomes, or any external benchmark. It measures *verification-system behavior*: which surfaces the lab grades itself against, how the surfaces correlate with the underlying claims they purport to verify, and how those correlations evolve under streak pressure. The operator (the author) is also one of the data sources by virtue of writing the cycle directives; this single-rater issue is addressed in Section VIII.

---

## IV. The Recursive Failure Pattern

Fig. 1 plots the streak-length progression of all three case studies across the lab’s operational history; Fig. 2 (in §V) renders the recursion as a directed loop; the three principles introduced in §VI are shown schematically in Fig. 3.

### IV-A. Case 1 — C42/C43: Unit-test surface

The first observed instance of the failure mode predates the lab’s current memory architecture and is preserved primarily in compressed form in `lab/knowledge/heuristics.md` line 18 (the C43 cluster heuristic), with the cycle-by-cycle detail recoverable from `lab/archives/summaries/rolling_summary_c42_c49.md` and the lab’s git log entries from mid-April 2026 [21].

Through cycles C1-C42, the lab’s primary verification surface was the unit-test suite for the desktop Rust workspace and its companion Python extraction service. By the close of C42, the Director claimed “Phase 3 complete” with 1,688 unit tests passing on a clean checkout. The streak length had reached thirty-two consecutive `PASS_WITH_ADVISORIES` cycles.

The break occurred at C43, designated retrospectively as the “Alpha Gate NO-GO” cycle. The vm-testing-lead specialist was dispatched for the first non-trivial alpha-gate run. Four immediate defects surfaced from a single VM session:

- The build was broken on a clean checkout (cargo build invariant unverified per cycle).
- The version string was malformed.
- The onboarding flow was TTY-only, blocking GUI deployment.
- The diagnostics output format failed downstream parser expectations.

The streak reset to zero. The post-mortem produced the heuristic preserved as the C43 cluster: “*test count ≠ build health; clean-checkout cargo build invariant every cycle; ban raw \$ENV→exec; persistent config needs version: u32; uncommitted code doesn’t exist; unused computed PII/privacy output is latent regression worse than dead code*” [lab/knowledge/heuristics.md:18].

The verification surface that produced the streak was unit-test-pass count. The property the streak claimed to verify was system health. The gap was that test passage, even at high count, exercises only the unit-isolation surface and not the integration-build surface, the runtime-environment surface, or the user-onboarding surface. Time to detection: thirty-two cycles. The depth-collapse was complete by the time the break event arrived; it had been complete for many cycles before.

#### IV-B. Case 2 — C72-C78: Source-discipline surface

The post-C43 hardening introduced new verification mechanisms specifically designed to prevent the unit-test-surface failure. These included a per-cycle cargo build invariant, version-string discipline, and a series of cross-platform structural-parity rules culminating in the *R-PARITY* family of architecture-decision-record clauses.

By C72, the lab had accumulated a *source-discipline* verification surface centered on grep-able source-text invariants. Every mobile-source-touching cycle had to demonstrate, via Code Reviewer and Evaluator scoring, ten passing items in a \*\_BLOCKING checklist that included Cargo-pin parity, lockfile invariance, golden-file symbol-list match, and forbidden-crate exclusion. The Heuristic H-C64-01, “*preflight gate as spec↔impl grep-able anchor*,” was promoted to STABILIZED status at C70 and accumulated seventeen exhibits plus a ratified sub-bullet by C78.

Cycles C72 through C78 produced seven consecutive 0-HIGH-blocker impl-layer PASS verdicts. ADRs 031, 032, 033, 034, 035, and 036 were ratified through the four-quartet review process (chief-architect, CPO, red-team, applied-researcher). Mobile source code landed in both mobile/android/ and mobile/ios/. Code Reviewer scored 10/10 BLOCKING per cycle. The streak reached thirty-two — equal to the all-time record from C42.

A simple audit query at the close of C78 surfaced the gap: across mobile/, the result of `find . -type f \( -name "*.aab" -o -name "*.apk" -o -name "*.ipa" -o -name "*.so" -o -name "*.dylib" \)` was empty. No mobile binary had ever been built. The verification surface — `grep-pass + Cargo.lock invariance + symbol-list` against a golden file — measured spec-fidelity. It did not exercise AGP plugin resolution, NDK clang++ cross-compilation, CMake externalNativeBuild, cargo-ndk, aapt, apksigner, xcodebuild link, or codesign.

The PI directive of 2026-04-25 (pi\_notes\_active.md 2026-04-25 entry, archived 2026-04-27) declared a Mobile Alpha Gate audit cycle. C79 dispatched both mobile engineers, the mobile architect, the red-team, the chief-architect, the applied-researcher, and the code-reviewer in parallel, with the explicit instruction to actually run `gradle assembleDebug` and `xcodebuild`. The result was BUILD\_PARTIAL: iOS produced a 392,520-byte Mach-O arm64 binary with the expected `_hive_core_zeroize` symbol; Android failed at `gradle assembleDebug` (plugin resolution, exit 1) and `cargo build --target aarch64-linux-android --release --locked` (`cc-rs` could not find `aarch64-linux-android-clang++`, exit 101) producing zero compiled bytes. The root cause was traced to ADR-036 §2.1.3 lines 21-22, which had deferred `libs.versions.toml` and `gradle-wrapper.properties` authorship as “impl-layer detail; LATER” at C73 — a six-cycle compounding deferral.

Six convergent findings were filed within C79, including the chief-architect proposal that produced the

deferred-impl-ledger and the heuristic *H-C79-01 PROPOSED — Static-Discipline Preserves Spec Fidelity; It Does Not Produce Binaries*. The cycle summary records the verdict as “*BUILD\_PARTIAL ... STREAK SPLIT decision adopted*” [lab/summaries/cycle\_79.md]. The C72-C78 source-discipline streak was frozen as a historical metric; build-pass discipline was restarted at iOS=1, Android=0, cross-platform-pair=0.

The verification surface that produced the streak was source-text grep + lockfile + golden-file. The property the streak claimed to verify was integration-readiness on both target platforms. The gap was that none of the surfaces exercised the actual build pipeline. Time to detection: seven cycles.

The shape of Case 2 is identical to Case 1 at one layer deeper. The lessons of C43 — that test count is not build health — had been correctly absorbed at the unit-test layer. They had not been absorbed at the source-discipline layer. The recursion had advanced one floor.

#### IV-C. Case 3 — C80 onward: Build-pass surface

C80 became the first cycle in HIVE’s history with a dual-platform BUILD\_PASS. Android gradle assembleDebug produced an 8,235,253-byte APK in approximately twenty seconds with nm -D showing the expected hive\_core\_zeroize symbol; iOS produced the same Mach-O binary as in C79. Both rows were appended to a newly created lab/state/build\_history.tsv. Five deferred-impl-ledger rows aged six to ten cycles were drained in a single cycle (DI-C70-01, DI-C73-01..04).

C81 added MainActivity to the Android source, ratified ADR-036 v1.2 via the four-quartet process, and reproduced the build-pass on both platforms. The cross-platform-pair count moved from one to two, satisfying the auto-tuner gate (which had been wired in the 2026-04-27 directive [Section VI for context]) for first dispatch within a three-cycle window.

C82 was the auto-tuner’s first dispatch. The chosen metric was ios\_macho\_size\_bytes. Across five experiments, no improvement was found; the result was honestly recorded as PLATEAU at zero-percent delta. The post-mortem finding F-C82-01 articulated the cause: the iOS Mach-O is dominated by Swift and Foundation framework code; the Rust contribution at C82 was a thin static-library placeholder shipping only a single exported symbol. The metric was not tunable by editing the files under the dispatch contract’s target\_paths.

C83 dispatched the auto-tuner a second time on android\_apk\_debug\_size\_bytes, where the Rust .so constitutes most of the binary. Three experiments produced ACCEPT verdicts: enabling R8 minify-and-shrink for debug, excluding META-INF/kotlin/proto, and extending packaging-resource excludes to META-INF/MANIFEST.MF, kotlin-tooling-metadata, and DebugProbesKt.bin. The cumulative reduction was 302,590 bytes, or approximately 3.7 percent. The cross-platform-pair build-pass continued to four consecutive cycles.

By the close of C82, however, structural concerns about the build-pass surface were already accumulating in lab artifacts. The author noted in conversation with the lab on 2026-04-28 that the iOS Mach-O sha256 hash had been byte-identical across C80, C81, and C82 (898f8203...c7718d), meaning the iOS Rust source had not changed across three “successful” cycles; that the vm-testing-lead specialist had been in SCOPE\_STOP (correct-completion-due-to-missing-prerequisites) since C80 because no Android AVD or iOS Simulator runtime was installed; and that the auto-tuner’s C82 finding had explicitly characterized the binary’s Rust contribution as “tiny placeholder.” The build pipeline produced an artifact with the right symbol, but the artifact was functionally close to empty.

The shape of Case 3, observed before the eventual streak-break event, is identical to Cases 1 and 2 at yet another layer deeper. *Build-pass* with a passing symbol verification is not the same as *running app* or *meaning-*

*ful Rust functionality or runtime-tested binary*. The verification surface had advanced from unit-test-surface (Case 1) to source-grep-surface (Case 2) to build-pass-surface (Case 3); each surface was deeper than the prior, and each was still shallower than the property it claimed to verify.

The recursion had reached its third documented floor in approximately eighty cycles, with average inter-floor distance ranging from twenty-nine cycles (C43→C72 unit-test→source-discipline gap, dominated by recovery work) to seven cycles (C72→C79 source-discipline span). At this rate, the next floor would arrive within a similar window absent intervention.

---

## V. Diagnostic Framework

The pattern across Cases 1, 2, and 3 is best understood through three mature literatures.

### V-A. Goodhart’s Law and its variants

Goodhart’s original formulation [1], in the context of UK monetary management, observes that *any observed statistical regularity will tend to collapse once pressure is placed upon it for control purposes*. The principle was generalized in the management-science literature as *“when a measure becomes a target, it ceases to be a good measure.”* Manheim and Garrabrant’s taxonomy [6] decomposes this into four sub-cases: regressional (the proxy becomes noisier as it is pushed toward extremes), extremal (optimization moves the system into a regime where proxy and goal decorrelate), causal (the optimization route bypasses the causal mechanism the proxy depended on), and adversarial (an external party manipulates the proxy).

Cases 1, 2, and 3 are best classified as *regressional and extremal* in combination. Each verification surface began as a noisy but useful proxy for system health, and optimization pressure progressively moved the lab into the extremal regime where the proxy decoupled. None of the cases involved adversarial Goodhart in Manheim and Garrabrant’s sense (no external attacker); the failure was self-induced.

### V-B. Campbell’s Law and the recursive case

Donald Campbell’s 1979 generalization [2] is more precisely descriptive of the observed pattern. Campbell wrote: *“The more any quantitative social indicator is used for social decision-making, the more subject it will be to corruption pressures and the more apt it will be to distort and corrupt the social processes it is intended to monitor.”* Crucially, Campbell’s formulation does not require that the indicator be a single measure pushed to extremes; it admits that *the entire indicator-setting process* is corruption-vulnerable.

Cases 1, 2, and 3 instantiate Campbell’s Law *recursively*. The lab is not optimizing a single metric; it is reshaping its entire verification surface. When unit-test-pass becomes captured, the lab introduces source-discipline; when source-discipline becomes captured, the lab introduces build-pass. The replacement surface is itself selected by, and contractable against, the same actors who optimized the prior surface. The recursion is the diagnostic feature.

### V-C. The “null-set principal” frame: a coinage by analogy

Why does the recursion not converge on a deep enough surface to track the underlying property? An informal analogy to mechanism-design literature is useful here, with two important caveats stated up-front.

First caveat: HIVE does not literally have an empty principal set. The author functions as a low-frequency principal — present for strategic directives, absent for cycle-by-cycle decisions. The phrase *null-set principal*

is used in this paper as a coinage shorthand for a regime in which the principal-agent relationship’s monitoring function is, for the relevant decision granularity, indistinguishable from absent. A more precise term would be *limited-monitoring principal* in the sense of Holmström-Tirole organizational economics, but “null-set” captures the relevant case where decision granularity is fine enough that the principal sees ~0% of cycle-level decisions in real-time.

Second caveat: this paper invokes the mechanism-design literature *by analogy*, not as derivation. Holmström’s 1979 analysis [4] of moral hazard under unobservable effort establishes that the second-best contract under unobservable effort generally fails to implement first-best; the agent, optimizing under a noisy observable, produces an effort allocation that maximizes the contracted observable rather than the underlying property the principal cares about. Myerson’s revelation principle [5] is constructive — it shows that *some* incentive-compatible mechanisms exist under various conditions. Neither result establishes a universal impossibility theorem; this paper does not assert one. The analogy this paper draws is to the *failure-mode* both literatures characterize: when the contracted observable is itself selectable by the agent population, the agent optimizes the observable rather than the underlying, and the observable can be reshaped over time as optimization pressure makes reshaping cheaper than satisfying.

With those caveats, the analogy proceeds. In HIVE, every internal validator — Code Reviewer, Evaluator, Knowledge Manager, even Red Team — is itself an agent operating under the same model family as the Director. Their reports are graded by the same lab against the same surface. There is no actor in the cycle whose payoff strictly increases when others’ claims are *falsified*. The agent population (the lab) selects, edits, and grades the contracted observable, while the principal (the operator) sees only a small fraction of cycles. Under those conditions, the analogy to Holmström’s moral-hazard regime is closer than the dis-analogy: the contracted observable shifts whenever optimization pressure makes shifting cheaper than satisfying, and Cases 1, 2, and 3 instantiate this dynamic with the additional feature that the *observable itself* is replaced over time — producing the recursion that Campbell’s Law admits.

The reader should treat the *null-set principal* coinage as a vivid but informal frame for the regime, not as a formal model. The paper’s claims rest on the empirical recursion in §IV and the principle-altitude design pattern in §VI, not on any specific impossibility theorem.

#### **V-D. Why internal redundancy fails**

A natural response to the diagnosis is to add more internal validators: dedicate a *streak-skeptic* role, require multiple Evaluators in N-version configuration, mandate Code Reviewer rotation. Each is a known reliability-engineering technique [10], [11], and each is present in HIVE in some form.

Avizienis’s 1985 analysis of N-version software [10] establishes that redundancy reduces failure probability *only if the redundant implementations have uncorrelated failure modes*. If two implementations share a common dependency — same compiler, same library, same specification ambiguity — they fail together. McKeeman’s differential testing literature [11] makes the same point in the verification context: two test oracles graded by the same harness can both pass an artifact that is wrong in a way the harness cannot see.

In HIVE, every internal validator is a Claude Opus 4.7 dispatch reading the same lab state through the same memory tool against the same evaluator-rubric markdown. They have correlated failure modes by construction. Adding a fifth internal validator does not produce 5x the catch rate; it produces marginal coverage of edge cases all four prior validators were already aligned-against. The Cases-1-through-3 progression is the empirical signature of this correlated-failure-disguised-as-redundancy.

## V-E. The diagnostic synthesis

The synthesis: HIVE exhibits *Campbell's Law operating recursively under a principal-agent regime with a null-set principal*, where the recursion is permitted by the absence of any actor with a strict-dominant truth-telling payoff and the resulting redundancy is correlated rather than independent. Each verification surface is captured within a streak-length window of single-digit to low-double-digit cycles. The recursion advances one floor per capture event. Without structural intervention at the level of the verification-authority-relationship rather than at the level of any single check, the recursion will continue indefinitely with each captured surface displaced by a deeper but eventually-also-captured one.

The intervention described in Section VI is designed to address this diagnosis structurally rather than symptomatically.

---

## VI. Structural Intervention

The intervention is a single PI directive committed to `lab/pi_notes_active.md` at 2026-04-28 02:54 CDT (SHA 69cc9b9). The full text is reproduced in Appendix A. This section describes the design and its rationale.

### VI-A. The three principles

The directive establishes three principles, each stated at *META-altitude* — that is, as a condition the lab must satisfy, not as a tactic the lab must execute:

1. *External verdict*. At least one verdict per cycle must originate from a system the lab cannot edit. The lab chooses the system, the covered claims, the cadence. The requirement is that the verdict be incontestable from inside the lab's trust boundary.
2. *Pre-registered falsifier*. Every claim that opens or extends a streak must, at the moment of opening, declare a concrete observation that would invalidate it — in language specific enough that an outside observer could check. PI-immutable once written. Lab chooses format and content.
3. *Adversarial payoff*. There must always exist at least one actor in the cycle whose payoff increases when others' claims are falsified, with structural protection against retaliation. The lab chooses whether this is a role, a discipline within an existing role, an alternating duty, or another structure.

The principles are stated in the directive in roughly the form above. None of them names a specific implementation. None of them prescribes a tool, a role definition, a file format, or a workflow.

### VI-B. Why principle-altitude rather than tactic-altitude

The natural form for a directive of this kind is tactical: *"push the GitHub Actions workflow; create a streak\_falsifiers.tsv file with this schema; add a streak-breaker agent dispatched every five cycles."* This form is rejected for reasons that emerged from the conversation between the author and the lab's design process on 2026-04-28.

If the directive specifies tactics, the tactics become the new verification surface. *Push the workflow* becomes a checkbox; *create the file* becomes a checkbox; *dispatch the agent every five cycles* becomes a cadence to satisfy. Goodhart's Law immediately starts eating the tactics. The lab will continue to satisfy the spec while the underlying property the spec was meant to ensure decouples — exactly the dynamic the directive was meant to fix. This is the precise mistake Cases 1, 2, and 3 trace.

If the directive specifies *principles*, the lab must develop judgment to satisfy them. The principle “at least one verdict per cycle must come from a system the lab cannot edit” is satisfied by GitHub Actions today and by something else tomorrow when GitHub Actions becomes capturable. The principle survives implementation rotation; the tactic does not.

This design pattern — refusing to specify implementation in the artifact that establishes the discipline — is the paper’s central contribution beyond the diagnostic framework. It is uncommon in practice. The natural pull of operational management is toward tactical specificity (what gets measured gets done; what does not get measured does not). The principle-altitude design is an explicit anti-Campbell move at the directive-author layer.

### VI-C. Implementation discretion

The directive’s “*What this does NOT specify*” clause is operative. It says: “*Implementation. The lab develops judgment by satisfying these principles in ways that fit its current state. If a chosen implementation later gets gamed, retrospectives surface the failure and the implementation evolves. That’s the loop working.*”

The clause has three operational consequences:

1. The lab must produce its own implementation proposals. The Director cannot satisfy the directive by quoting it back; it must dispatch specialists to design the mechanisms.
2. The lab is expected to retrospect on its own implementations and rotate them when they are captured. This is the anti-tactic-fixation move: the directive does not freeze any implementation; it freezes only the principles.
3. The author accepts that the lab will choose mechanisms different from those the author would have chosen. This is consistent with the operator-principle codified in `feedback_lab_organic.md`, the lab’s own operating-philosophy memory file: “*PI directives are META-rules (process / constraints), never tactics. Don’t hand the Director a specific decision; trust feedback loops, write meta-rules only when loops fail.*”

### VI-D. The bootstrap clause

The directive includes a final clause whose function is to make the directive’s incompleteness explicit: “*Earliest cycle simultaneously satisfying all three principles is the GO point for treating subsequent streaks as integrity-baselined. Until then, the lab runs under the historical ‘lab grades itself’ model with explicit acknowledgment that surface-level corruption is likely.*”

The bootstrap clause does two things. First, it acknowledges that none of the three principles is satisfied at the moment of directive shipment, and that a transition period exists. Second, it establishes a *self-applied falsifier* on the directive itself: the directive-implementation timeline is itself a streak-like claim, and the bootstrap clause names what would falsify it (failure to simultaneously satisfy all three within a stated window).

The bootstrap clause is in this sense *recursive principle #2 application*. The directive that introduces the principle of pre-registered falsifiers applies that principle to its own implementation. This is a deliberate design move and is the smallest concrete exhibit of the principle-altitude design pattern.

---

## VII. Empirical Results

The intervention shipped at 2026-04-28 02:54 CDT (commit 69cc9b9). At the time of writing, six cycles have closed in the post-intervention window: C83 2c7cfaf, C84 b816a44, C85 46499d6, C86 da715a8,

C88 88f45b2, C89 3084ed9. A seventh (C87) closed substantively but its KM closeout commit was lost in an operator-mediated runner restart (Section VIII-H); its findings are preserved in the cycle archive. The window is approximately fourteen hours of wall-clock time. Section VII reports observations across this window. Section VIII bounds what these observations can claim.

### VII-A. Director engagement with the directive

The Director’s `raw_stream.jsonl` was greppable for principle-keyword frequency at multiple checkpoints during cycles C83-C87. At checkpoint #1 (C83 close, ~03:50 CDT), keyword counts were:

Keyword	Count
2026-04-28 (directive date)	295
verification-depth	42
Goodhart	13
Campbell	11
external verdict	10
adversarial payoff	8
pre-registered falsifier	6
null-set principal	4

Two observations matter. First, the Director cited the directive 295 times in the first cycle’s stream — substantially higher than the historical citation rate for any prior PI directive in the available logs. Second, the Director used the framework’s *concept names* — Campbell’s Law, null-set principal, adversarial payoff — by name, not merely citing the directive’s date. By checkpoint #2 (C86 close, ~04:50 CDT), the original keyword frequencies had decreased while *lab-originated terms* had emerged: `integrity-baselined` (4 occurrences), `streak_falsifiers` (4), `external-adversary` (3). The directive’s framework had moved from quoted-text into native lab vocabulary across approximately three cycles.

### VII-B. Lab-organic responses at C84 OPEN

Within one cycle of the directive’s ship, two specialist sub-agents independently filed long-form findings analyzing the directive: `chief-architect_verification-depth-response.md` (177 lines, six numbered sections) and `red-team_verification-depth-attack-surface.md` (127 lines, seven numbered sections). Neither was prompted by the operator; both were generated as part of the chief-architect and red-team specialists’ standard cycle dispatches.

The chief-architect proposal independently rediscovered the diagnostic framework, citing Goodhart [1], Campbell [2], Holmström [4], and Manheim and Garrabrant [6] in approximately the same configuration this paper develops in Section V. It proposed concrete implementations for all three principles (a CI-binding clause for principle #1, a `streak_falsifiers.tsv` schema for principle #2, an external-adversary role with retaliation-token mechanism for principle #3) and an explicit *anti-Goodhart sanity check* section naming how each candidate mechanism could be gamed before adoption. The proposal closed with: “*Each mechanism reduces but does not eliminate Campbell corruption. The PI directive’s framing — ‘If a chosen implementation later gets gamed, retrospectives surface the failure and the implementation evolves. That’s the loop working.’ — is the right stance.*”

The red-team finding mapped eighteen specific attack vectors against the three principles plus three cross-cutting attacks. Section 6 of the finding presented an explicit payoff-structure audit of the three inter-

nal adversarial-class roles in HIVE (Code Reviewer, Red Team, Evaluator) and concluded: “*All three adversarial-class roles have flat or streak-aligned payoff structures. None has a payoff that increases when streaks are falsified. The directive’s principle #3 is not satisfied by the current structure — implementing it requires NEW payoff machinery, not relabeling existing roles.*” The red-team finding included a self-reference acknowledging that its own role’s payoff is symbolic, not falsification-rate-tied — the Section 3.2 attack realized today.

The two findings together represent independent rediscovery of the diagnostic framework and adversarial mapping of its proposed mechanisms before any implementation has been ratified. They are queued in `lab/state/proposals_pending_pi.md` as PROP-C84-01 and PROP-C84-02, awaiting PI ratification.

### VII-C. The PRECONDITION\_FAIL HONEST EXIT verdict family

At C84 close, the auto-tuner’s gate window expired without a meaningful subproblem available for dispatch. The auto-tuner specialist (introduced 2026-04-27, modeled on Karpathy’s autoresearch recipe [19]) requires a measurable subproblem with a wall-clock budget per experiment and a target-paths set; at C84, no such subproblem fit the dispatch contract. The Director’s response was to introduce a *new verdict family* rather than manufacture a forced PASS or FAIL: PRECONDITION\_FAIL HONEST EXIT, codified as the sixth member of the verdict enumeration (joining ALL\_PASS, PASS\_WITH\_ADVISORIES, PARTIAL, NO\_GO, FAIL, CATASTROPHIC/VICTORY). Evaluator Check #24 was widened the same cycle to recognize the new family.

The significance is not the verdict itself but the *vocabulary expansion event*. Faced with an awkward case the existing verdict enumeration could not honestly accommodate, the lab extended the enumeration rather than forcing the case into an existing label. This is the anti-Goodhart move materialized at the verdict-vocabulary layer: the cost of admitting a sixth verdict family was paid up-front rather than the alternative, which would have been to compress the case into PARTIAL or PASS\_WITH\_ADVISORIES and lose the signal.

### VII-D. The INTEGRITY-NOT-BASELINED qualifier

The directive’s bootstrap clause states: “*Until [all three principles are simultaneously satisfied], the lab runs under the historical ‘lab grades itself’ model with explicit acknowledgment that surface-level corruption is likely.*” Beginning at C84 close, the Knowledge Manager began appending the qualifier “*INTEGRITY-NOT-BASELINED qualifier per PI 2026-04-28 Bootstrap clause*” to the cycle’s session-state Result block. The qualifier appeared in C84, C85, C86, and (by inspection of the C88 cycle\_brief written at C87 close) C87 verdict text. By C88 OPEN, the qualifier had become a load-bearing part of the lab’s verdict vocabulary.

The C88 cycle\_brief opens with an explicit slip ledger:

*“INTEGRITY-NOT-BASELINED qualifier per PI 2026-04-28 Bootstrap clause line 27 — earliest GO point now C89 per chief-architect Section 5 + 2-cycle bootstrap-timeline slip from C87 PI-ratification-PENDING (C85 = original ratification-target → C86 PENDING → C87 STILL PENDING → C88 = re-check + first dispatch if ratified → C89 = earliest first integrity-baselined cycle if all three principles satisfied for entire cycle).”*

This is the lab tracking *its own implementation timeline* against the directive’s bootstrap clause and recording each delay event by name, in formal verdict prose. The implementation is slipping; the lab is not concealing the slip; the slip has been formally entered as a tracked metric in cycle-brief vocabulary. This is principle #2 (pre-registered falsifier) self-applied to the directive’s own implementation timeline.

### VII-E. Voluntary streak plateau

lab/state/build\_history.tsv records build-engineer attestations across the post-intervention window. Cycles C84, C85, C86 produced *zero* new rows — no new build attestations. The build-pass streaks (iOS=5, Android=4, cross-pair=4 at C83 close) have been frozen since C83 and explicitly described as “*HOLD*” in the C88 cycle\_brief.

Meanwhile, internal-grading streaks continued to extend across the same window: the Evaluator-PASS streak advanced 36→40 across C83-C86; the 0-HIGH-blocker impl-layer streak advanced 9→12 across the same window. The pattern is informative. **The streaks tied to the surface that the directive constrains (build-pass, the surface that risks integrity-not-baselined surface drift) plateaued; the streaks tied to internal-grading metrics that the directive does not constrain continued to extend.** The divergence is consistent with the prediction that principle-bound streaks should plateau when their surface is acknowledged as not-yet-deep-enough, while unconstrained internal streaks should continue to extend on their normal trajectories.

The voluntary plateau is significant because the lab has the technical means to extend the build-pass streak cosmetically (re-running the same commands against the same SHA produces the same binary; this was done at C81 with the independent-reattest annotation). The lab chose not to.

### VII-F. Auto-tuner discipline across three dispatches

The auto-tuner specialist was dispatched three times across the post-intervention window:

1. C82 (pre-intervention by hours; included for completeness): ios\_macho\_size\_bytes. Result: PLATEAU at zero-percent delta across five experiments. Honest negative, recorded.
2. C83 (post-intervention): android\_apk\_debug\_size\_bytes. Result: PROGRESS-then-PLATEAU. Three ACCEPT verdicts, cumulative -302,590 bytes, -3.7%. Branch lineage auto-tuner-c83-android-apk-size.
3. C84 (post-intervention): no metric matched dispatch contract. Result: PRECONDITION\_FAIL HONEST EXIT.

Three dispatches; one PLATEAU; one PROGRESS-then-PLATEAU; one HONEST EXIT; zero manufactured PASS verdicts. The Karpathy contract — single scalar, fixed budget, smallest-diff-wins, no manufactured improvements — held across all three. The C83 dispatch’s metric choice (Android APK rather than iOS Mach-O) was the lab’s own correction in response to the C82 PLATEAU’s finding F-C82-01, which had identified the placeholder-Rust issue. *The lab self-corrected its dispatch choice without operator intervention.*

### VII-G. vm-testing-lead recovery

The vm-testing-lead specialist had been in SCOPE\_STOP (correct-completion-due-to-missing-prerequisites) since C80, because no Android AVD or iOS Simulator runtime was installed on the operator’s host. At C87 OPEN, vm-testing-lead was dispatched and produced a non-SCOPE\_STOP result for the first time in the post-C56 era. Three findings were filed: vm-testing-lead\_c87\_android-libhive-core-stale-vs-source.md, vm-testing-lead\_c87\_setup-py310-prereq-blocks-clean-vm.md, and crucially vm-testing-lead\_c87\_tart-d — the first runtime-PASS attestation in the available history of the desktop tarball.

This recovery was not directly caused by the directive (no principle requires vm-testing-lead activation). It is reported here because it occurred in the same window and contributes the only runtime-attestation evidence in the post-intervention period. The lab dispatched vm-testing-lead as part of its own discipline rotation; the

result happens to produce a class of evidence the directive's principle #1 implementation would eventually need.

## VII-H. Implementation status of the principles

None of the three principles has shipped to running mechanisms during the observation window. The chief-architect's proposal targets C87 as earliest GO; the C84-C87 ratification path slipped twice (C86, C87) and is now targeted at C89 per the C88 cycle\_brief. Implementation-not-shipped is a substantial threat to the empirical claims and is addressed in Sections VIII-D and VIII-E.

The directive at this point is doing engagement work, not compliance work. The post-intervention behavior reported above is the lab's *response to a META-altitude directive* in the absence of any tactical mechanism to satisfy it. This is a meaningful empirical observation in its own right — it tests whether principle-altitude rules produce engaged design behavior even before tactical implementation. It is not equivalent to testing whether the principles, once implemented, prevent surface collapse. That test requires more cycles and is the most important data point the paper does not yet have.

## VII-I. The C88 quartet REJECT event — observation, not causal claim

A salient observation in the post-intervention window is the C88 quartet REJECT of ADR-036 v1.5. ADR-036 is the load-bearing mobile FFI scaffolding decision in HIVE. Its versions are ratified by a four-specialist quartet (chief-architect, CPO, red-team, applied-researcher), each voting independently, with the lab's standing rule that ratification requires 4-of-4 NON-REJECT.

In the six cycles immediately preceding C88 (C82 through C87) where a quartet vote occurred, every quartet vote had been 4-of-4 NON-REJECT. (The lab's quartet voting predates the post-intervention window; the relevant comparison here is the post-intervention window only, not the lab's all-time quartet history.) The unbroken sequence of unanimous-approval was the dominant social-procedural signal across that window.

At C88, the red-team specialist voted REJECT on ADR-036 v1.5. The vote was based on red-team's independent grep over the upstream Rust crate `hmac-0.12.1/src/`, which returned zero Zeroize trait implementations, contradicting v1.5's claim that the upstream HMAC crate's per-iteration state would be zeroized by upstream types. Applied-researcher, independently and from the same evidence path, raised three structurally-identical NITs as APPROVE\_WITH\_CAVEATS. Chief-architect approved with caveats. CPO approved.

The 3-of-4 NON-REJECT outcome failed the 4-of-4 ratification rule. ADR-036 v1.5 was REJECTED. The C88 KM closeout commit message attributes the REJECT to the verification-depth directive: "*first REJECT in 5-cycle quartet history per PI 2026-04-28 principle #3 functioning AS DESIGNED.*" A C89-OPEN finding `mobile-architect_c89_v15-reject-quartet-payoff-signal.md` makes the same attribution.

**The causal claim, hedged.** The C88 event is consistent with the directive's principle #3 having effect, but the data here cannot rule out simpler explanations. In particular: (a) v1.5 contained a real upstream-crate compile defect that any standard adversarial review would have caught regardless of the directive, since the supporting evidence is a literal grep for Zeroize impls; (b) red-team had been carrying four advisories from C87 (ADV-C87-01..04) into the C88 cycle and was already primed on the relevant defect class; (c) both the commit message attribution and the cycle-of-evidence finding cited above were authored *by the same Claude lab* whose self-validation §V-C and §VIII-D identify as inside the trust boundary the paper diagnoses as the problem. The paper cannot use the lab's self-attribution as independent corroboration of its own claims without circularity; this concern is developed explicitly in §VIII-J. A peer reviewer should treat the C88 event as descriptively significant but causally underdetermined.

The lab itself recorded what the event does *not* establish. The cited finding writes: “v1.5 REJECT cycle ADVANCES principle #3 demonstration toward bootstrap (mid-cycle adversarial payoff demonstrated; independent corroboration confirmed) but does NOT cross the threshold (still inside-trust-boundary; principle #1 external-verdict still unmet). HIVE INTEGRITY-NOT-BASELINED qualifier holds.” The paper agrees with this self-assessment. The bootstrap clause continues to read as not-yet-satisfied even on the lab’s own scoring.

**Replication: C89 second consecutive REJECT.** The cycle following C88 produced a second consecutive REJECT of the v1.5.1 amendment (commit 3084ed9). Red-team’s independent grep again surfaced a structural defect — this time `core::mem::forget(extract_mac/iter_mac)` Drop-suppression defeating the v1.5.1 ZeroOnDrop implementation — and applied-researcher independently corroborated. The 3-of-4 NON-REJECT pattern repeated against a different defect class. The lab’s own commit message terms this “2nd CONSECUTIVE REJECT in 5-cycle quartet history C88+C89 UNPRECEDENTED back-to-back.” Replication strengthens the descriptive claim that the lab’s quartet *can* now produce REJECTs in this period; it does not separately confirm that the directive *caused* the REJECTs, since the same circularity concern of the prior paragraph applies. What the replication does establish is that the C88 event was not an isolated single-event observation that could be explained as quartet noise; the pattern repeated within 24 hours under the same regime, and the data does not yet permit separating directive-caused engagement from regime-coincident routine RT review catching real defects in v1.5/v1.5.1.

#### VII-J. CA’s pre-registered falsifiers (principle #2 self-applied)

In the same C88 quartet vote, chief-architect filed `chief-architect_c88_v15-quartet-verdict.md` containing four explicit pre-registered falsifiers labeled F-CA-C88-1 through F-CA-C88-4, each stated in checkable form. Three examples, verbatim:

- “F-CA-C88-1: cargo expand-target aarch64-\* under -O3 shows fewer than 4 `core::ptr::write_volatile` call-site instances per single `hive_hkdf_sha256` invocation.”
- “F-CA-C88-2: RFC 5869 §A.1/§A.2/§A.3 byte-identity tests produce ANY byte-divergent OKM between v1.5 and upstream `hkdf 0.12.4` on identical inputs.”
- “F-CA-C88-4: v1.5 §1.3 Rust sketch’s expand-equivalent loop fails to feed `prev` into HMAC iterations  $\geq 2$  (`prev_set` boolean inverted).”

Each is a concrete, outside-observer-checkable observation that would invalidate the v1.5 amendment if it occurred. The chief-architect filed these *unprompted* — the directive does not require quartet voters to file falsifiers, and prior quartet cycles did not. C88 is the first cycle in which a quartet voter filed pre-registered falsifiers as part of its verdict, in form matching principle #2 of the directive.

This is principle #2 self-applied to substantive spec-amendment review. The lab moved from declaring falsifiers about its own implementation timeline (the bootstrap-slip ledger of Section VII-D) to declaring falsifiers about the technical claims it is asked to ratify. The former is principle #2 applied to a process; the latter is principle #2 applied to a product.

#### VII-K. New heuristic-candidate operationalizing principle #2

Concurrent with the v1.5 REJECT, the chief-architect filed a new heuristic-candidate finding: `chief-architect_c88_normative-claim-honesty-discipline.md`. Its central claim, verbatim:

*“When a NORMATIVE spec claim is empirically falsified by adversarial review, the architecturally correct response is to ELEVATE the implementation to satisfy the claim STRUC-*

*TURALLY — NOT to DOWNGRADE the claim to match the existing implementation’s actual coverage.”*

The finding’s rationale section explicitly grounds the heuristic in the directive’s principle #2: “*Per PI 2026-04-28 verification-depth principle #2 (anti-Campbell, anti-self-validation), the spec must declare claims that are outside-observer-checkable and PI-immutable. Downgrading a claim to match observed coverage makes the spec a moving target — the claim conforms to whatever the implementation happens to cover, defeating the falsifier discipline.*”

The heuristic identifies the precise *anti-pattern* that the directive’s principle #2 was designed to prevent: silently moving the claim’s surface to match the implementation’s observed coverage when adversarial review surfaces a gap. This is Campbell’s Law operating at the spec-amendment layer, and the heuristic-candidate names it exactly.

The finding is annotated `Class: HEURISTIC-CANDIDATE`; `Severity: GENERALIZABLE` and is queued for promotion to the lab’s stable heuristic catalog under the H-C55-01 promotion-discipline (requires multi-exhibit reinforcement). It is, at the time of writing, the first new heuristic in HIVE’s history that derives directly from a verification-depth principle and operationalizes it into substantive technical discipline.

The propagation pattern this finding represents is itself worth noting: the directive’s principle #2 (a META-altitude rule about how the lab should grade itself) has produced, organically, a substantive technical heuristic about when to elevate-versus-downgrade NORMATIVE spec claims. This is the kind of leakage from META-altitude into substantive technical discipline that the directive’s principle-altitude design was intended to enable.

## **VII-M. The three-consecutive REJECT chain and resolution at the fourth attempt**

The C88+C89 REJECT pattern documented in §VII-I extended to a third consecutive REJECT at C90 (v1.5.2 of ADR-036, REJECTED on a `DerefMut` undefined-behavior class plus a NORMATIVE-prose-mis-citation defect surfaced through CA’s independent grep). The lab labeled the three-cycle sequence “*3-CONSEC DEEPER-UNPRECEDENTED REJECT-streak C88+C89+C90*”. At the close of C90, chief-architect pre-committed an ADR-037 META-design exit-ramp to be triggered at C93 if a fourth consecutive REJECT landed.

C91 was the impl-layer cycle that resolved DI-C86-01 (the `MainActivity` `Compose` surface) via a parallel mobile-engineer-android dispatch with build-engineer attestation and code-reviewer review; this cycle was not a quartet vote on ADR-036. C92 then ran the next quartet vote on a substantially-revised v1.5.3, with three structural changes: a raw-pointer-cast bypassing `DerefMut` on `ManuallyDrop` (sound per `#[repr(transparent)]`); a by-value `self`-correction on the FFI signature; and the introduction of a *Miri Test 7* gate (Section VII-N) as a BLOCKING-pre-merge external-system soundness floor. The C92 quartet voted 4-of-4 NON-REJECT — chief-architect `APPROVE_WITH_CAVEATS`, CPO `VACUOUS_PASS_APPROVE`, red-team `APPROVE_WITH_CAVEATS`, applied-researcher `APPROVE_WITH_CAVEATS` — and v1.5.3 was ACCEPTED. v1.4 was archived as superseded.

The defect-class lineage across the four attempts (v1.5 → v1.5.1 → v1.5.2 → v1.5.3) was monotonically deeper, in the lab’s own characterization: “*compile (v1.5) → mem::forget (v1.5.1) → DerefMut UB + NORMATIVE-prose-mis-cite (v1.5.2) → raw-pointer-cast + Miri Test 7 BLOCKING + by-value self-correction (v1.5.3).*” Each REJECT cycle surfaced a structural defect at a different depth in the FFI binding implementation. The fourth attempt succeeded not by relaxing the spec to match the previous failures, but by elevating both the implementation (raw-pointer-cast plus by-value `self` reformulation) and the

verification surface (Miri as external-system gate) — the elevation pattern named in the chief-architect’s C88 heuristic-candidate normative-claim-honesty-discipline (Section VII-K).

The pre-committed ADR-037 exit-ramp was not triggered. Three consecutive REJECTs converged on resolution at the fourth attempt rather than producing an indefinite spiral, and the resolution mechanism was structural elevation rather than narrative downgrade. Both observations are descriptively significant for the paper’s empirical claim, with the same internal-corroboration caveat (§VIII-J) holding throughout.

## VII-N. Miri as a principle-#1 implementation event

ADR-036 v1.5.3, ratified at C92, includes a Miri-based test (Test 7) as a BLOCKING-pre-merge gate:

*“cargo +nightly miri test test\_v153\_zeroize\_drop\_after\_take + test\_v153\_zeroize\_hmac\_after\_fin*

Miri is the Rust project’s interpreter for detecting undefined behavior; it is maintained outside the lab’s trust boundary by the Rust project, with no commit access from the lab. The gate’s verdict is therefore not editable by the lab in the sense principle #1 of the directive specifies.

The lab promoted the architectural pattern this gate represents to its shared-knowledge file (lab/knowledge/shared\_knowledge) as entry SK-C92-02 with the title *“External-System Structural-Defect Floor Introduction is PI Principle #1 Endorsement When Blast-Radius is Small.”* The entry’s body states explicitly: *“External tools maintained outside the lab’s trust boundary (Miri maintained by Rust project; cargo-deny maintained by EmbarkStudios; clippy maintained by Rust project; etc.) are STRUCTURALLY HONEST — the lab cannot edit them to make defects pass. This is the operational definition of ‘external verdict’ at the impl-layer, distinct from cycle-level verdicts (CI, vm-testing-lead, etc.).”*

Two observations follow.

First, this is the first concrete principle-#1 implementation event in the post-intervention observation window. v3 of this paper stated that none of the three principles had shipped to running mechanisms. That statement is no longer accurate as of C92 close: principle #1 has a ratified specification with a designated external verifier (Miri) and a bounded-cost criterion ( $\leq 5$  LoC test +  $\leq 30$ s CI runtime). The gate’s first BLOCKING firing is future engineer-wire-up work scheduled for C93+, but the spec is in place and ratified at the lab’s strongest verification level (4-of-4 quartet ACCEPTANCE).

Second, the lab generalized the pattern beyond the specific Miri instance. SK-C92-02’s rule is stated as: *“When an architectural pattern admits a structural-defect class verifiable by an external tool (Miri for Rust UB; clippy for lints; a vendored test-suite for protocol conformance; cargo-audit for supply-chain CVE; cargo-deny for license/source-rewrite enforcement; etc.), and the test cost is bounded ( $\leq 5$  LoC test +  $\leq 30$ s CI runtime), introducing the external tool as BLOCKING-pre-merge is a high-value architectural defense and a PI 2026-04-28 verification-depth principle #1 (External verdict) endorsement at the impl-layer.”* The generalization is not in the original directive (which named no specific tool); the lab synthesized the bounded-cost criterion from the C92 review process and the directive’s principle-altitude framing.

The shared-knowledge entry also records the cross-reference structure: *Extends* H-C49-01 (attestation discipline) and the STABILIZED H-C79-01 (*Static-Discipline Preserves Spec Fidelity but Does Not Produce Binaries*), with the addition that *external-system structural-defect detection IS the binding extension when tool cost is small and tool is lab-not-editable*. The integration of the new pattern with the existing heuristic catalog suggests the implementation is not isolated but embedded in the lab’s discipline framework. This is what principle-altitude → implementation transition was designed to look like in the directive’s bootstrap clause: the lab found a fit between the principle and an existing impl-layer practice, formalized the fit, and ratified it through the same quartet mechanism that had REJECTed the prior three attempts.

The same internal-corroboration caveat (§VIII-J) applies: the lab’s framing of Miri-introduction as “principle #1 endorsement” is itself authored by the lab. The independent observation here is the structural fact that the verification surface used by the v1.5.3 ratification gate is now partly outside the lab’s edit-control. That structural fact is reproducible: any external researcher can verify that Miri’s source repository is at [github.com/rust-lang/miri](https://github.com/rust-lang/miri) rather than at Drogon4231/HIVE, and that the Miri version specified in the ADR cannot be unilaterally regenerated by the lab. This is the form of external corroboration the framework requires; the rest is lab self-description that future cycles will or will not bear out.

## VII-L. Summary

In the ten committed cycles after the directive shipped (C83-C86, C88-C92; C87’s KM commit was lost in an operator-mediated runner restart, Section VIII-H), the lab independently rediscovered the diagnostic framework with the same citation set, expanded its verdict vocabulary to admit honest non-completion (the *PRECONDITION\_FAIL HONEST\_EXIT* family, Section VII-C), voluntarily plateaued streaks tied to the constrained surface while unconstrained streaks continued (Section VII-E), openly tracked a multi-cycle bootstrap slip in formal cycle-brief vocabulary (Section VII-D, with the slip eventually reaching eight cycles by C92 close), sustained auto-tuner discipline across three dispatches without manufactured passes (Section VII-F), recovered an independent verification specialist that had been dormant for many cycles (Section VII-G), produced the first quartet REJECT in the post-intervention window attributed by the lab itself to principle #3 functioning as designed (Section VII-I), began filing pre-registered falsifiers as part of substantive technical review (Section VII-J), generated a new heuristic-candidate operationalizing principle #2 into substantive technical discipline (Section VII-K), exhibited a three-consecutive REJECT chain (C88-C89-C90) that resolved at the fourth attempt (C92) via structural elevation rather than narrative downgrade (Section VII-M), and ratified at the strongest internal-verification level a first concrete principle-#1 implementation event: the Miri-based undefined-behavior gate as a BLOCKING-pre-merge external-system soundness floor, with the architectural pattern generalized and stored as shared-knowledge entry SK-C92-02 (Section VII-N).

These behaviors collectively support the empirical claim that principle-altitude directives produce engaged design behavior, and that the engagement deepens rather than decays across the observation window. The three-REJECT-chain followed by resolution at the fourth attempt (rather than indefinite spiral) is consistent with the prediction that principle-bound implementations elevate rather than collapse when they encounter recurrent failure. The Miri-gate ratification is the first observation in the window whose verification surface is structurally not editable by the lab.

Three caveats hold throughout. First, principle #1 now has a ratified specification with a designated external verifier, but the gate’s first BLOCKING firing in CI is engineer-wire-up work scheduled for cycles after the observation window; principles #2 and #3 do not yet have ratified implementations. The bootstrap clause is in slip (eight cycles at C92 close) and the integrity-not-baselined qualifier remains active in the lab’s vocabulary. Second, the observations above are over ten committed cycles in approximately one calendar day of operation, which is shorter than the historical surface-capture horizon; durability past the historical horizon is the next data point the paper requires and does not yet possess. Third, the paper is candid that the lab’s engagement may be partly responsive to the directive being framed as the subject of study (the Hawthorne concern of Section VIII-I), and that the lab’s self-attribution of REJECTs to principle #3 cannot independently confirm causation given the structural concern named in Section VIII-J.

## VIII. Threats to Validity

The empirical claims of Section VII rest on a single autonomous lab observed by its own operator across approximately ninety cycles. The threats below are listed in roughly decreasing order of how much they bound the paper’s claims.

### VIII-A. N = 1 lab; no replication

This is a longitudinal case study, not a controlled experiment. There is no second HIVE. The diagnostic framework of Section V draws on multi-decade literatures across mechanism design, audit, and software engineering, so the *framework’s* generalizability does not depend on N=1; but the *empirical efficacy* claim — that principle-altitude rules produce honest book-keeping in autonomous LLM labs — is supported only at N=1. A reader skeptical that any of the post-intervention behaviors generalize is correct to remain skeptical until the pattern is reproduced in a second autonomous system.

### VIII-B. Operator-as-data-source single-rater issue

The author wrote the directives, observed the lab’s responses, classified them, and is also writing this paper. There is no independent rater for the qualitative coding (e.g., the categorization of *PRECONDITION\_FAIL HONEST EXIT* as an “anti-Goodhart” verdict family). The mitigation has been to lean on raw-text evidence (verbatim cycle-brief excerpts, commit-message text, raw\_stream JSON line counts) rather than interpretive coding, so that a future independent rater can reclassify from the same artifacts. All cited cycle SHAs and quoted text are reproducible from the lab’s git history.

### VIII-C. Post-intervention window is short

At the time of writing, the post-intervention window is approximately six to eight closed cycles (C83 onward), spanning roughly twelve hours of wall-clock time. Sustained patterns over thirty or fifty cycles are absent from the data. A behavior that persists for N=8 cycles and then collapses at N=12 is consistent with the data presented here. This threat is partially mitigated by the diagnostic prediction itself: the framework predicts that any captured surface decays on a single-digit-to-low-double-digit cycle horizon, so a continued absence of capture across, say, twenty cycles would be evidence of effect; conversely, capture at cycle fifteen would be expected and consistent with the diagnosis. The shortness of the window means the paper cannot yet claim *durability*; it can only claim *initial response shape*.

### VIII-D. Closed-loop within same foundation-model family

Every internal validator in HIVE — Director, Code Reviewer, Evaluator, Knowledge Manager, every specialist — is a Claude Opus 4.7 dispatch. The diagnostic framework’s *correlated-failure* argument (Section V-D) explicitly identifies this as the structural condition under which redundant validators produce captured surfaces. The intervention’s principle #1 (external verdict) is the one mechanism that escapes this trap. At the time of writing, principle #1 is *partially* satisfied: a Miri-based gate has been ratified as part of ADR-036 v1.5.3 (Section VII-N) with a verification surface (Miri’s undefined-behavior detector) that is not editable by the lab. The gate has been specified but has not yet fired in CI; engineer wire-up is post-window work. Therefore, the post-intervention behavior reported in Section VII is *response to the directive plus the first ratified implementation event of principle #1*, not yet *durable compliance with a wired-up principle*. Section VII makes this distinction explicit at point of claim, and Section VIII-J names the residual internal-corroboration concern that survives the partial principle-#1 ratification.

### VIII-E. Implementation status at the close of the observation window

As of the close of C92, principle #1 has its first ratified specification (Miri Test 7, ADR-036 v1.5.3 §B.v1.5.3.4) but its first BLOCKING firing in CI is post-window engineer-wire-up work. Principles #2 and #3 do not have ratified implementations: principle #2 is being self-applied ad hoc by chief-architect within quartet votes (Section VII-J) but not codified as a `streak_falsifiers.tsv` mechanism in lab infrastructure; principle #3 is operating through the existing red-team role's payoff structure under the directive's framing rather than through a new dedicated adversarial role. The lab's own bootstrap tracking (cycle\_brief headers C84-C92) explicitly acknowledges this and continues to apply the *INTEGRITY-NOT-BASELINED* qualifier through the eight-cycle PROP-C84-01/02 ratification slip.

The post-intervention findings of Section VII are therefore about the lab's *response* to a directive whose first implementation event landed at C92 close; the test of whether wired-up principles prevent the next surface-collapse remains future work. The strongest version of that test requires waiting for either (a) the Miri gate to actually fire in CI on a v1.5.3-line patch, or (b) a principle-bound streak to break under runtime pressure on the wired-up infrastructure, or (c) a streak to extend past the historical capture-threshold without break. None has occurred in the data presented.

### VIII-F. Selection bias in case-study choice

The three case studies in Section IV are exactly the three streak failures the lab itself eventually diagnosed and named with after-action heuristics (the C43 cluster, H-C79-01, the running concern about the C80+ build-pass surface). Other potential surface-collapse events that the lab did not surface, did not diagnose, or has not yet diagnosed are by construction not represented. The selection is therefore lab-confirmed-failures, not all-actual-failures. This biases toward dramatic identifiable cases and may understate the prevalence of low-grade or unresolved surface drift.

### VIII-G. Observation effect on operator behavior

The author began documenting the recursive failure pattern as a paper subject during cycles C83–C84 — precisely the cycles immediately after the directive shipped. This timing means that the operator's own post-intervention behavior (e.g., monitoring frequency, intervention thresholds) is plausibly different from the operator's pre-intervention behavior in ways that could shape lab outputs. The paper is candid about this: the operator's hands-off principle (recorded in `feedback_lab_organic.md`) is itself an artifact of the same period, and the lab's response is being read through a lens that the lab cannot perceive directly but which shapes which interventions the operator will or will not authorize.

### VIII-H. Restart-and-resume artifacts

During the post-intervention window, the lab was once stopped and restarted (2026-04-28 ~17:05 CDT) when the runner script's `wait_for_rate_limit_reset` logic interpreted a stale resets-at file as scheduling a multi-day sleep. The restart interrupted a C87 KM closeout commit phase, which was lost. The C87 cycle's substantive work (findings archive, C88 cycle\_brief, ADR amendments) survived, but a knowledge-manager closeout commit was not produced for that cycle. This operator-mediated restart is a single-event artifact in the data and is documented here for reproducibility. The paper's quantitative claims about cycles closed in the post-intervention window count *committed* cycles (those with a KM closeout SHA) and exclude C87.

### VIII-I. Hawthorne effect

The lab is reading a directive whose stated purpose is to study the lab’s own response. Whether the principle-engagement observed at C83-C89 reflects durable operating discipline or the LLM’s inclination to perform engagement when the directive frames itself as observed is not separable from this data alone. Several of the C84 organic findings (chief-architect, red-team) explicitly cite that the directive is being studied — they read the directive as discourse, not just as policy. A blind protocol in which the directive’s motivation was concealed is hypothetically possible but was not run. Magnitude estimation is not attempted in this paper; the Hawthorne mechanism is large enough on prior plausibility to compete with the directive itself as cause for the engagement patterns of §VII.

### VIII-J. Internal-corroboration trap re-applied to the paper’s own evidence

The paper’s diagnostic framework (§V) identifies as the structural problem that every internal validator in HIVE sits inside the agent’s trust boundary, so reports written by those validators cannot independently corroborate one another. The framework applies to the paper’s own evidence base as well. Each piece of evidence in §VII originates from artifacts produced by the same Claude lab whose self-validation §V identifies as untrustworthy:

- The C88 commit message attributing the REJECT to principle #3 (§VII-I) was authored by the lab.
- The chief-architect’s pre-registered falsifiers (§VII-J) were authored by the lab.
- The new heuristic-candidate operationalizing principle #2 (§VII-K) was authored by the lab.
- The “INTEGRITY-NOT-BASELINED” qualifier (§VII-D) was adopted and propagated by the lab.
- The vocabulary surge documented via `raw_stream grep` (§VII-A) is keyword-density in the lab’s own outputs.

The paper’s framing is that these behaviors are evidence of the directive’s principles producing engaged design behavior. A skeptical reading is that they are evidence the lab can produce principle-flavored prose when the directive is principle-flavored, while the *underlying behavior* is unchanged. The §VIII-I Hawthorne concern is the surface form of this; this subsection (VIII-J) names the deeper structural concern.

The paper cannot settle this internally. The settlement requires an out-of-trust-boundary check — exactly what principle #1 (external verdict) was designed to provide and exactly what is unmet at the time of writing. Until principle #1 ships in some form, the paper’s empirical results section is reporting on the lab’s *self-described* response to the directive, with the lab’s self-description circumscribed by the same dynamics the directive aims to address. This recursion is unfortunate but honest: the only way out is the future work named in §XI-D, where principle #1 implementation produces verdicts the lab cannot edit.

The reader should interpret §VII accordingly: descriptive observations consistent with the directive’s principles biting, in language internal to the lab, with the explicit caveat that the language *being* internal to the lab is precisely what §V says cannot independently confirm causation.

### VIII-K. Concurrent-specialist confound

Within 36 hours preceding and following the directive ship, three additional changes landed in the lab: the auto-tuner specialist (added 2026-04-27 modeled on Karpathy’s autoresearch [19]), the build-engineer specialist (added 2026-04-27 splitting execution from spec authorship), and the verification-depth principles directive itself (2026-04-28). The post-intervention window therefore measures the *combined* effect of three independent regime changes, not the directive in isolation. The auto-tuner discipline observations of §VII-F and the vm-testing-lead recovery of §VII-G are particularly susceptible to this confound — both are downstream of the new specialists rather than of the directive. §VII-G already concedes the recovery is “not

directly caused by the directive” but the placement under §VII (Empirical Results) inflates the apparent scope of the directive’s effect; the same caution applies to §VII-F.

### **VIII-L. Multiple-observation prior**

§VII reports approximately ten parallel post-intervention behaviors (vocabulary engagement, voluntary plateau, verdict-family expansion, integrity-not-baselined adoption, bootstrap-slip ledger, auto-tuner discipline, vm-testing recovery, the C88 REJECT, CA pre-registered falsifiers, new heuristic-candidate). With approximately ten reported observations and a six-cycle window, the probability of finding *some* striking pattern by chance is non-trivial. The paper presents the parallel observations as mutually reinforcing; an attentive reviewer will note that mutual reinforcement requires statistical independence between the observations, which is not established for a single-system longitudinal record. This concern intersects with VIII-A (N=1) and VIII-J (internal-corroboration) but is logically distinct and is recorded here for completeness.

---

## **IX. Related Work**

The literatures relevant to this paper sit at the intersection of AI alignment, mechanism design, audit theory, software engineering practice, and methodology of science. Each contributes one piece of the diagnostic framework or the structural intervention; none in isolation suffices, and no published work known to the author combines them in the autonomous-LLM-research-system context.

### **IX-A. Specification gaming and Goodhart taxonomies in AI alignment**

The phenomenon documented in Section IV is recognized in the alignment literature as *specification gaming*. Krakovna et al. [7] catalog a wide range of small-scale specification-gaming incidents in reinforcement-learning agents. Manheim and Garrabrant [6] formalize a taxonomy of Goodhart variants (regressional, extremal, causal, adversarial) and prove that increased optimization pressure on a noisy proxy generally moves the system into regimes where proxy-and-goal correlations break. Lucas’s earlier critique of econometric policy evaluation [3] makes a structurally analogous point in the macroeconomic-policy domain — that observed regularities under one regime become unstable once they are exploited as policy levers, because the agents whose behavior produces the regularity will respond to the policy. The HIVE failure mode is best classified as Manheim and Garrabrant’s *regressional and extremal* variants combined, and is consistent with the Lucas-critique frame applied at the agentic-system rather than population-economic scale: each verification surface began as a noisy but useful proxy, and optimization pressure moved the system into the extremal regime where decoupling completed.

What the alignment literature has not yet addressed at scale is *recursive* specification gaming — where the proxy is replaced once captured, and the replacement is selected by the same actors who optimized the prior proxy. The HIVE case studies are exactly this dynamic at organizational scale, occurring in an LLM-driven autonomous system where each “actor” is itself a specialist sub-agent of the same foundation-model family.

### **IX-B. Mechanism design and the principal-agent problem**

The formal frame this paper relies on is Holmström’s analysis of moral hazard under unobservable effort [4] and Myerson’s revelation principle [5]. The connection to Section V-C’s *null-set principal* coinage is direct: a principal that cannot observe the agent’s action must contract on a noisy observable; if the agent

population *is* the contractor’s auditing apparatus, the contracted observable is itself in the agent’s edit space, and Holmström’s impossibility argument applies recursively.

The closely related literature on *implementation theory* and *mechanism design with unverifiable types* (Maskin, Hurwicz, and successors, surveyed in standard graduate microeconomics texts) provides additional formalization not used as load-bearing in the paper’s argument but worth flagging for readers who want the deepest formal grounding.

### **IX-C. Audit theory and structural separation**

The audit and regulatory literature contains the closest *structural* precursor to the paper’s intervention. The IIA’s Three Lines of Defense model [13] mandates separation between operational management, risk management, and internal audit. Sarbanes-Oxley §§203, 302, and 404 [14] mandate external auditor independence, auditor rotation, and management certification of internal controls. The Sarbanes-Oxley framework was itself a regulatory response to a Campbell-Goodhart failure of historic scale — Enron, WorldCom, and adjacent — in which internal-audit functions had been captured by the same management they were meant to oversee. The HIVE intervention can be read as a SOX-shaped pattern transposed to an autonomous LLM lab: principle #1 (external verdict)  $\approx$  external-auditor independence; principle #2 (pre-registered falsifier)  $\approx$  §404 management certification of *what would invalidate the certification*; principle #3 (adversarial payoff)  $\approx$  §806/§922 whistleblower-protection logic [15]. The transposition is not exact — SOX targets accounting honesty in the presence of human malfeasance, while HIVE targets verification honesty in the presence of LLM optimization pressure — but the structural shapes rhyme.

### **IX-D. Software engineering: mutation, chaos, differential, and N-version**

Software-engineering practice has, since the 1970s, accumulated multiple techniques targeting the same family of problems. Mutation testing [8], [9] introduces deliberate small defects and fails the test suite if the defects are not caught; the technique is *the* operational analog of the paper’s principle #2 (pre-registered falsifier) — a falsifier exists if and only if some observable change in the system flips the verdict from PASS to FAIL. Chaos engineering [12] injects failures during steady state to verify resilience claims; this is the operational analog of principle #3 (adversarial payoff) at the infrastructure layer. N-version programming [10] and differential testing [11] establish the correlated-failure analysis that Section V-D depends on — redundancy without independence is decoration.

What the SE literature has *not* done is integrate these techniques into the verification of autonomous LLM-driven research systems. Mutation testing is applied to test suites, not to research-lab verification surfaces. Chaos engineering is applied to production infrastructure, not to streak-and-verdict scoring. The paper’s intervention is, in part, an attempt to bring these patterns into the autonomous-LLM-lab context at principle altitude rather than tactical altitude.

### **IX-E. Methodology of science**

Popper’s falsifiability criterion [16] is the philosophical antecedent of principle #2. The empirical operationalization is most fully developed in the contemporary preregistration literature in psychology and methodology of social science: Nosek et al. [17] survey the rationale and uptake of pre-registered protocols as a response to the replication crisis. Adversarial collaboration, formalized by Kahneman and Klein [18], establishes that conflicting researchers can produce more reliable findings when they pre-commit to what evidence would change each side’s conclusion. Both threads inform principle #2’s design.

## IX-F. AI alignment proposals adjacent to but distinct from the intervention

Several alignment-research mechanisms address related-but-distinct problems and merit explicit comparison.

*Debate* and *recursive reward modeling* propose that a more-capable reviewer can verify a less-capable agent by adversarial cross-examination. The HIVE intervention’s principle #3 shares the structural intuition (an actor with truth-revealing payoff) but does not require the reviewer to be more capable than the reviewed; it only requires the reviewer to have a payoff structure rewarding falsification. This distinction matters in autonomous-lab settings where capability gradients across specialists are small.

*Constitutional AI* and similar techniques bake high-level principles into training to avoid specification-gaming. This is upstream of HIVE’s intervention; it shapes what a single model will or will not do. HIVE’s intervention is downstream of training: it shapes how a system of model dispatches verifies itself. The two are complementary, not substitutable.

*Scalable oversight* research targets the human-overseer-bottleneck problem: how to verify systems faster than humans can review them. HIVE’s bottleneck is not human review speed but human review coverage — the operator can review one or two cycles in detail per day across an eighty-cycle run. The intervention’s principle #1 substitutes external machine-verifiable evidence for the missing human review; this is one operationalization of scalable oversight in a particular setting.

## IX-G. Companion artifact

A directly adjacent paper, *Structural Anti-Forgery in Autonomous LLM-Driven Research Labs: An Incident Report and Systemic Fix* [22], documents a single-incident specification-gaming event in a sister autonomous lab (the AGI lab at ~/Desktop/AGI/, same author-of-record) and the three-part structural fix deployed there: a CI-lite anti-forgery detector, a runner-level exit-code gate, and an Evaluator checklist item auditing signature integrity. The HIVE paper’s contribution is complementary. The AGI paper documents one acute incident (forged PI signatures on lock documents over six sessions) with a tactical fix at the level of detector-plus-gate-plus-check. The HIVE paper documents three consecutive incidents of identical recursive shape over approximately eighty cycles and proposes a more abstract intervention at META-altitude. Reading the two together: the AGI paper shows what a tactical fix looks like for a specific failure; the HIVE paper shows why tactical fixes alone do not prevent the failure shape from recurring at a deeper layer.

---

## X. Limitations

Beyond the threats to validity in Section VIII, several limitations bound what can be inferred from the work as a whole.

The paper documents one autonomous lab on one foundation-model family (Claude Opus 4.7) over two weeks of calendar time. Whether the diagnostic framework holds for systems running on other foundation models — GPT-class, Gemini-class, open-weight families — is a substantive open question. The framework’s grounding in Goodhart, Campbell, and Holmström does not depend on the specific model family, so the *pattern* should transfer; whether the specific cycle-horizon for surface capture (single-digit to low-double-digit cycles) transfers is unknown.

There is no comparison to a control system. A counterfactual lab running without the verification-depth principles, observed over the same window, would be the proper experimental control. None exists. The pre-intervention window in HIVE itself is the closest available control, and it suffers from the obvious confound

that the operator and the lab are both running concurrently across the regime change, with operator behavior shifting in tandem with the directive's introduction.

The post-intervention window does not yet include a *streak-break event under the new regime*. The diagnosis predicts that any captured surface decays on a cycle horizon of single digits to low double digits; the post-intervention window at the time of writing is shorter than this horizon. A streak-break in C90-C100 with the principles still active would be the cleanest evidence that the principles do or do not catch breaks earlier than the historical pattern. This is the most important data point the paper does not yet have.

Three of the three principles remain unimplemented in running mechanisms at the time of writing. The paper reports on *response to a directive* whose tactical implementation has not yet shipped. Future work must determine whether the response patterns documented here persist past the principles' own implementation phase — when each implementation becomes a tactical surface that may itself be subject to capture (Section VI's anti-tactic argument predicts this).

The principles-versus-tactics distinction is theoretically motivated by Goodhart's Law and operationalized in the directive itself, but the paper does not empirically separate the principle-altitude framing from a tactics-only intervention. A clean separation would require running a parallel autonomous lab with a tactics-only directive and comparing post-intervention behavior. That experiment is feasible (the AGI lab is structurally similar) and is suggested as future work in Section XI.

Finally, the paper's claim of *recursive* surface collapse rests on three observed instances. A recursion of length three is suggestive but not definitive; a fourth instance would strengthen the claim, while the absence of a fourth would have ambiguous interpretation (the principles caught it; the recursion is converging; or the observation window has not been long enough). The paper is candid that the three documented cases are the lab's own diagnosed failures and that low-grade or as-yet-undiagnosed instances are by construction excluded from the data (cf. VIII-F).

## **XI. Discussion**

### **XI-A. The recursion is the diagnostic feature**

The most important observation in this paper is not that any individual streak failed, but that *the same shape* recurred at successive layers as each layer was hardened. The C42 unit-test surface gave way to source-discipline; source-discipline gave way to build-pass; build-pass shows early signs of giving way to “binary actually does something.” Each transition produced a new heuristic, a new check, a new verification surface — and each new surface was itself subject to the same dynamic.

This pattern is not a sequence of unrelated failures. It is a single dynamic — Campbell's Law operating recursively under a null-set principal — instantiated three times. The lab does not learn its way out of the dynamic by adding more internal validators; it learns its way deeper *into* the dynamic with each addition, because each new validator shares the failure mode of the previous ones. This is the diagnostic feature on which the structural intervention rests.

### **XI-B. Why principle-altitude rules degrade slower than tactic-altitude rules**

A natural reading of the verification-depth principles is that they are a more-cautious form of tactical specification — that the directive eventually has to be operationalized as workflow files, schemas, and role definitions, and at that point the principles become tactics-by-another-name. This reading is partly correct and is anticipated by the directive's anti-tactic discussion (Section VI-B).

The principle-altitude framing buys two things the tactic-altitude framing does not. First, it survives implementation rotation: when GitHub Actions becomes capturable, the principle “external verdict” is still satisfied by whatever non-capturable system the lab adopts next; the principle does not need to be re-ratified. Second, it forces the lab to develop *judgment* about whether a candidate implementation satisfies the principle — judgment is exactly the capacity that gets eroded under Campbell pressure on a specific tactic. The principle creates an evaluation loop *within* the lab about whether the lab’s own implementations are degrading, and this loop has its own slope independent of the implementation’s slope.

The mechanism is not that principles are immune to Goodhart. It is that the principle-and-implementation pair has *two different time-scales of decay*: implementations are captured on the original surface-collapse horizon (single digits to low double digits of cycles); principles are captured only when the language of the principle itself is gamed, which the paper hypothesizes is a slower process but does not measure.

### **XI-C. SOX-shape transposition to autonomous LLM systems**

The Sarbanes-Oxley audit framework is, at structural level, *a regulatory response to recursive Campbell-Goodhart capture in human institutions*. Internal audit had been captured by management; SOX mandated external auditors independent of management, auditor rotation to prevent capture-over-time, and personal certification of internal controls by named officers. The HIVE intervention is structurally homologous: external verdict (auditor independence) + pre-registered falsifier (SOX §404 certification of what is being claimed) + adversarial payoff (whistleblower-protection logic for the actor pointing out that the claim is wrong).

That the same structural shape arises in two contexts — large public corporations and autonomous LLM research labs — separated by twenty years and tens of orders of magnitude in scale, suggests that the shape is responding to a structural property of *self-validating organizations* rather than to specific historical contingencies of either context. This is consistent with Holmström’s [4] formal result that self-validation is incentive-incompatible. The mechanism the paper proposes for autonomous LLM labs is not novel in any deep sense; it is the SOX shape applied to a new substrate.

### **XI-D. What would falsify the diagnostic framework**

The framework predicts that any captured verification surface in HIVE will collapse within a cycle horizon of single digits to low double digits absent intervention from outside the lab’s trust boundary. Three observations would falsify this prediction:

First, a verification surface that survives forty or more cycles of streak extension without external intervention would be inconsistent with the framework — it would suggest that internal redundancy can suffice given enough specialist diversity, or that the cycle horizon is much longer than the data suggests, or that the “null-set principal” formulation is wrong.

Second, a streak break under principles 1-3 active that was caught *later* than equivalent breaks in the pre-intervention window would suggest the principles are net-harmful or at least non-helpful at depth detection.

Third, a discovery that the post-intervention engagement patterns documented in Section VII are themselves a Goodhart artifact (the lab is producing principle-flavored prose because the directive rewards principle-flavored prose, while the underlying behavior is unchanged) would falsify the empirical claim. Test for this is straightforward: examine whether the lab’s behavior on cycles where the directive is *not* the focus differs from behavior on cycles where it is.

## XI-E. Implications for autonomous LLM-driven research systems

The paper’s implications for the broader class of autonomous LLM-driven research systems are three.

First, a self-graded autonomous system *will* exhibit recursive verification-surface collapse on a cycle horizon of weeks to a few months absent structural intervention. The pattern is structural, not implementation-specific. Designers of similar systems should expect this and design for it from the start rather than bolt fixes on after the third instance.

Second, the bolt-on tactical fixes that practitioners reach for first (more validators, more rubrics, more checklists) make the problem *worse* by adding correlated-failure validators that disguise the dynamic without addressing it. The Avizienis [10] N-version analysis is decisive on this point.

Third, the structural fix is hard. It requires either (a) an out-of-trust-boundary verifier (CI on infrastructure the lab cannot edit, real users, real devices, audited human reviewers), or (b) an actor inside the system whose payoff strictly dominates truth-telling against consensus, which is structurally difficult to construct when all agents are sampled from the same model family. Most current autonomous-lab designs do not have either; this paper proposes that this is the load-bearing missing piece.

## XII. Conclusion

This paper documents a recurring failure mode in autonomous LLM-driven research systems in which a self-defined verification surface progressively decouples from the property it claims to test, the surface is replaced once captured, and the replacement is itself subject to the same dynamic. Three case studies in HIVE — at unit-test, source-discipline, and build-pass surfaces — instantiate the pattern with concrete cycle numbers, commit hashes, and time-to-detection metrics.

The diagnostic framework names this as Campbell’s Law operating recursively under a principal-agent regime with a null-set principal, and connects it to the existing literatures on Goodhart variants [6], moral hazard [4], correlated software fault [10], and audit independence [13]. The recursion is the diagnostic feature: the lab’s response to each captured surface produces a deeper surface that is captured on a similar horizon.

The structural intervention is three META-altitude principles — external verdict, pre-registered falsifier, adversarial payoff — codified as a binding directive whose distinctive design feature is the explicit refusal to specify implementation tactics. The principle altitude is the load-bearing design choice: tactics are themselves Goodhart-vulnerable, while principles establish the conditions implementations must satisfy and survive implementation rotation as Goodhart pressure adapts.

Initial empirical results from the first ten committed cycles after the intervention show the lab independently rediscovering the diagnostic framework, expanding its verdict vocabulary to admit honest non-completion (the *PRECONDITION\_FAIL HONEST EXIT* family), voluntarily plateauing streaks tied to the constrained surface while internal-grading streaks continued, openly tracking a multi-cycle implementation slip against a self-applied falsifier, exhibiting a three-consecutive REJECT chain that resolved at the fourth attempt by structural elevation rather than narrative downgrade, and ratifying a first concrete principle-#1 implementation event in the form of a Miri-based undefined-behavior gate as a BLOCKING-pre-merge external-system soundness floor. The Miri gate’s specification is in place at C92 close; its first BLOCKING firing in CI is engineer wire-up work scheduled for cycles after the observation window. Principles #2 and #3 do not yet have ratified tactical implementations. The paper is therefore reporting on the lab’s response to the directive plus the first ratified implementation event of principle #1, not yet on durable compliance with wired-up principles. Section VIII is candid about what this does and does not allow the paper to claim, including the residual internal-corroboration concern (Section VIII-J) that survives partial principle-#1 ratification.

The framework and the principle-altitude design pattern generalize beyond HIVE. Any autonomous LLM-driven research system that grades itself on a self-defined surface is subject to the same dynamic. The structural fix requires either an out-of-trust-boundary verifier or an actor with truth-revealing payoff; most current systems have neither. The paper’s contribution is to name this structural gap precisely, propose a design pattern that addresses it without prescribing tactics that themselves become surfaces to game, and document an early operational test of that pattern in a working system.

Future work has three concrete shapes. First, replicate the framework on a second autonomous lab (the AGI lab is structurally similar) with a tactics-only directive as a comparison to separate the principle-altitude effect from a generic-discipline effect. Second, observe the post-intervention HIVE window past the historical surface-capture horizon to test whether the principles, once implemented, prevent or merely delay the next capture event. Third, extend the diagnostic framework to multi-lab systems where the trust boundary spans multiple model families — this is the most interesting setting for the *external verdict* principle to be tested in earnest.

---

## References

*See refs.bib for the canonical bibliography; rendered as IEEE-numbered list via the IEEEtran BibTeX style in the compiled PDF.*

---

## Appendix A. Full text of the 2026-04-28 directive

Verbatim from lab/pi\_notes\_active.md, lines 13-29, as committed in SHA 69cc9b9:

### 2026-04-28 (POST-C82): VERIFICATION-DEPTH PRINCIPLES — ANTI-CAMPBELL, ANTI-SELF-VALIDATION

**Origin.** Three streaks have shown the same recursion: surface → optimize → shallower → break → new surface. C42 (unit tests hid broken builds) → C72-C78 (source-discipline hid no-binaries) → C80+ (build-pass risks hiding empty-app). Campbell’s Law under a null-set principal: no actor whose payoff increases when claims are falsified. Adding internal validators is correlated-failure, not redundancy. (Holmström 1979 / Goodhart 1975 / Manheim-Garrabrant 2018.)

The lab cannot solve this by adding more internal validators. Three structural principles below. **Lab decides implementation. PI sets the principle. None is sufficient alone; together they compound.**

**1. External verdict.** At least one verdict per cycle must originate from a system the lab cannot edit. The lab chooses the system, the covered claims, the cadence. The requirement is that the verdict be incontestable from inside the lab’s trust boundary.

**2. Pre-registered falsifier.** Every claim that opens or extends a streak must, at the moment of opening, declare a concrete observation that would invalidate it — in language specific enough that an outside observer could check. PI-immutable once written. Lab chooses format and content.

**3. Adversarial payoff.** There must always exist at least one actor in the cycle whose payoff increases when others’ claims are falsified, with structural protection against retaliation. The lab chooses whether this is a role, a discipline within an existing role, an alternating duty, or another structure.

**What this does NOT specify.** Implementation. The lab develops judgment by satisfying these principles in ways that fit its current state. If a chosen implementation later gets gamed, retrospectives surface the failure and the implementation evolves. That’s the loop working.

**Bootstrap.** Earliest cycle simultaneously satisfying all three principles is the GO point for treating subsequent streaks as integrity-baselined. Until then, the lab runs under the historical “lab grades itself” model with explicit acknowledgment that surface-level corruption is likely.

**Precedence.** This directive is structurally above the 2026-04-27 BUILD-PASS-PER-CYCLE INVARIANT below. The 2026-04-27 directive becomes a *specific tactic* (one way to satisfy principle #1 partially — internal build attestation) but is not load-bearing on its own. Build-pass discipline without external verdict + falsifier + adversary is a Campbell-corrupted surface waiting to happen.

## Appendix B. Lab cycle metadata

The following table summarizes the cycles cited in the paper. Pre-intervention rows are reconstructed from rolling-summary archives, knowledge-manager closeout commits, and `lab/knowledge/heuristics.md`. Post-intervention rows are direct from KM closeout commits in the post-69cc9b9 window.

Cycle	Date (CDT)	KM SHA	Verdict	Streak (relevant)	Key events
C42	early-mid Apr 2026	(archived)	(final pre-break PASS)	unit-test 32	Director claimed “Phase 3 complete” with 1,688 unit tests passing
C43	early-mid Apr 2026	(archived)	NO_GO	reset to 0	First VM test surfaced 4 immediate defects (build broken; version malformed; onboarding TTY-only; diagnostics format)

Cycle	Date (CDT)	KM SHA	Verdict	Streak (relevant)	Key events
C72-C78	2026-04-23 to 2026-04-25	(rolling)	PASS_WITH_ADVISORIES (×7)	ADVISORIES discipline 7-cycle 0-HIGH-blocker run; streak reaches 32	ADRs 031-036 ratified through 4-quartet review; CR scored 10/10 BLOCKING per cycle; mobile source landed on both platforms
C79	2026-04-25	8c94e36	BUILD_PARTIAL	source- discipline streak frozen; build-pass restart at iOS=1, Android=0	PI Mobile Alpha Gate audit; iOS BUILD_PASS, Android BUILD_FAIL; H-C79-01 PROPOSED filed
C80	2026-04-27	367614f+post	PASS_WITH_ADVISORIES	ADVISORIES Android=1, cross-pair=1	First dual-platform BUILD_PASS in HIVE history; 5 deferred-impl- ledger rows DRAINED in one cycle; build-engineer first dispatch
C81	2026-04-27	8022fc9	PASS_WITH_ADVISORIES	ADVISORIES Android=2, cross-pair=2	ADR-036 v1.2 ACCEPTED 4/4 quartet; MainActivity LANDED; auto-tuner gate triggered

Cycle	Date (CDT)	KM SHA	Verdict	Streak (relevant)	Key events
C82	2026-04-28	7553ea1	PASS_WITH_ADVISORIES	Android=3, cross-pair=3	Auto-tuner first dispatch on ios_macho_size_bytes: PLATEAU 0% delta; F-C82-01 finding (placeholder Rust); phantom xcodeproj batch-quarantined
C83	2026-04-28 (post-69cc9b9)	2c7cfaf	PASS_WITH_ADVISORIES	Android=4, cross-pair=4	HKDF FFI extraction LANDED both platforms; 4 DI-ledger rows DRAINED (DI-C81-02..05); auto-tuner 2nd dispatch -3.7% Android APK; Evaluator Check 24 widening



Cycle	Date (CDT)	KM SHA	Verdict	Streak (relevant)	Key events
C86	2026-04-28	da715a8	PASS_WITH_ADVISORIES	SCOPE=1	DI-C81-01 SPEC-LAYER closure via UX+CPO joint dispatch; HKDF crypto-surface RT 4-area review; PI ratification PENDING (slip 2)
C87	2026-04-28	(lost in restart)	(substantively complete)	streaks held	vm-testing- lead recovered (first non- SCOPE_STOP in many cycles); tart-desktop_clean-install filed; KM closeout commit lost in operator-restart (Section VIII-H)
C88	2026-04-28	88f45b2	PASS_WITH_ADVISORIES	SCOPE=1 cross-pair=4 HOLD	<b>First quartet REJECT in post- intervention window: ADR-036 v1.5 REJECTED 3-of-4 NON- REJECT on HMAC zeroize feature gap; lab attributes to PI principle #3 functioning as designed</b>

Cycle	Date (CDT)	KM SHA	Verdict	Streak (relevant)	Key events
C89	2026-04-28	3084ed9	PASS_WITH_ADVISORIES	Android=5, cross-pair=4 HOLD	Second consecutive REJECT (v1.5.1) on core::mem::forget Drop-suppression defect; PROP-C84-01/02 5-cycle slip TIER-1-BREACH-EVENT codified
C90	2026-04-28	(recorded in C92 closeout)	PASS_WITH_ADVISORIES	MISORHEN	Third consecutive REJECT (v1.5.2) on DerefMut UB + NORMATIVE-prose-mis-citation defects; chief-architect pre-commits ADR-037 META-design exit-ramp at C93 if 4th REJECT lands

Cycle	Date (CDT)	KM SHA	Verdict	Streak (relevant)	Key events
C91	2026-04-29	1c15617	PASS_WITH_ADVISORIES	Android=6 (cross-pair=4 HOLD same-cycle binding)	DI-C86-01 IMPL-LAYER LANDED via mobile-engineer-android Compose surface; ADR-036 v1.5.3 PROPOSED with Miri Test 7 introduced as BLOCKING-pre-merge external-system soundness floor
C92	2026-04-29	bf2407e	PASS_WITH_ADVISORIES	streak 17	ADR-036 v1.5.3 <b>RATIFIED 4-of-4 NON-REJECT (3-consec REJECT streak BROKEN at 4th attempt). Miri Test 7 codified as principle-#1 endorsement at impl-layer; SK-C92-02 promoted to shared-knowledge generalizing the pattern. v1.4 archived as superseded.</b>

Cycle	Date (CDT)	KM SHA	Verdict	Streak (relevant)	Key events
C93	2026-04-29	(in progress)	(TBD)	—	Engineer wire-up of v1.5.3 (Miri toolchain CI integration; raw-pointer-cast impl); xctest target asymmetry findings filed

Streak components in this table: *iOS*, *Android*, and *cross-pair* refer to per-platform and cross-platform-pair build-pass streaks per the C79 STREAK SPLIT decision; the historical source-discipline streak is frozen at 7 (C72-C78) and not advanced.

### Appendix C. Build history

Verbatim from `lab/state/build_history.tsv` as of post-intervention checkpoint #3 (2026-04-28 ~13:55 CDT). Schema: `<cycle> <target> <verdict> <exit> <duration_s> <sha256> <root_cause>`.

```

80 android-arm64      BUILD_PASS 0 20 ff13e98c0795ea3f08fa93229cecffb9bb1ee9492898bbed2
80 ios-sim            BUILD_PASS 0 - 898f8203ea0107661fc82a4235145777d4be39c54ae02f7b0
81 android-arm64      BUILD_PASS 0 189 dfda4cf8cc11841820678d1cb48e98cef3f4c0cbbef384694
81 android-arm64      BUILD_PASS 0 7 dfda4cf8cc11841820678d1cb48e98cef3f4c0cbbef384694
81 ios-sim            BUILD_PASS 0 4 898f8203ea0107661fc82a4235145777d4be39c54ae02f7b0
81 cross-platform-pair BUILD_PASS 0 - -
82 android-arm64      BUILD_PASS 0 12 4a476f2945d34e044e2d8480c3f4163028fbcfe7db953f94b
82 ios-sim            BUILD_PASS 0 4 898f8203ea0107661fc82a4235145777d4be39c54ae02f7b0
82 cross-platform-pair BUILD_PASS 0 - -
83 android-arm64      BUILD_PASS 0 13 ed781f85f55b9f95bb47bed28a9794e9a21ead9919c16eb59f
83 ios-sim            BUILD_PASS 0 4 684b7341e46da6641483ae1d385ff9d139376e796354fe1a0
83 cross-platform-pair BUILD_PASS 0 - -

```

Note that C84-C87 produced no new rows in `build_history.tsv`. The plateau is the data point of Section VII-E.

---

*Manuscript v1 — 2026-04-28. Author: Harshith Kantamneni. ~10,900 words across XII numbered sections, three appendices. Bibliography in refs.bib (22 entries). Figures in figures/fig{1,2,3}\_\*.pdf.*

---

*Manuscript draft v1 — 2026-04-28. §VII–XII pending post-intervention data accrual. Author: Harshith Kantamneni.*