

Augmented Operational Decisioning with Ontology-Grounded Local Agents: A Fourteen-Run Empirical Study Against Real Construction-Operations Data

William White*¹

¹Bay West Labs | Little Bear Foundry Research

May 10, 2026

Abstract

This report consolidates fourteen experimental runs that empirically test whether a locally-hosted agentic substrate – consisting of TrustGraph, a Little Bear (LB) operational layer, a DoWhy/PyMC causal-reasoning sidecar, and LM Studio-hosted open-weights language models – can produce operationally-defensible decisions on real construction operations data drawn from a production OTTER system at Mac Construction. We adopted a two-database holdout discipline throughout: a frozen OLD-DB snapshot was used to construct the agent’s worldview, while a NEW-DB snapshot (never queried by the agent) was retained as ground truth. The program traverses four phases: *(i)* a diagnostic phase (runs 01–02) that surfaced canonical-prior dominance and a “substrate scope ceiling”; *(ii)* a position-paper and substrate-expansion phase (runs 03–05) that re-framed the product as an Augmented Operator Surface and committed to a *substrate-first* architecture; *(iii)* a verification and cadence phase (runs 06–10) that established the F2.x biases as substrate-bound rather than model-bound at $24 \times$ fixture scale; and *(iv)* a compounding-ontology and capability-primary phase (runs 11–14) that culminates in Hermes-4-70B-FP8 producing operationally-defensible foreman picks at a strict-equality rate of 66.7% and an operations-defensibility rate of 86.7%, with $\kappa = 0.95$ inter-rater reproducibility, zero hallucinations, and 100% compose reliability over 279 production-grade composes for \$8.50 total spend. Five of six foundational thesis claims are confirmed; the sixth (*cogent useful executables via autonomous tool use*) is reframed to *executable proposals via prompt-grounded RAG*. The bulk of this report (§13) is devoted to run 14, which carries the program’s most cofounder-defensible empirical signal.

1 Introduction

The Little Bear Foundry is a research program inside Bay West Labs that targets a single question: *can a local-first, ontology-grounded agentic system deliver decisions a working construction operations team would actually act on?* The answer matters because the alternative posture – cloud-hosted, autonomous, large-model agents that mutate production systems on their own – is operationally untenable for a regulated, schedule-driven, equipment-bound trade like utility-line construction.

*Correspondence: sheldon@baywestlabs.com. Affiliation: Bay West Labs | Little Bear Foundry.

The system under test combines four substrates: TrustGraph as the knowledge graph, Little Bear (LB) as the operational data and event-log layer, a DoWhy/PyMC reasoning sidecar for causal counterfactual queries, and LM Studio-hosted (and, for run 14, vLLM-hosted) open-weights models as the inference layer. Five production *role-agents* – foreman, scheduler, unblocker, form-classifier, and morning-brief – compose JSON proposals on this stack; proposals are gated by a Policy Engine before becoming executable mutations.

1.1 Thesis under test

The foundational thesis (originally articulated in the *Ontology-First Strategic Revision*, run-11 doc 19) asserts six joint claims about a locally-hosted ontology-grounded agentic system. Each subsequent run is scored against them:

1. **Locally hosted** – runs on operator-class hardware without API dependence in the hot path.
2. **Ontology-grounded** – decision-relevant facts are sourced from TrustGraph triples.
3. **Simulated real operations** – the fixture mirrors Mac Construction’s actual operational surface.
4. **Cogent, sensible, useful executable commands** – proposals are mutation envelopes an operator could execute directly.
5. **Plain-English redirectability** – the operator types free-form steering; the system extracts structure.
6. **Compounding ontology** – the redirect updates the ontology, and subsequent decisions read the update.

By run 14, five are confirmed; claim 4 is reframed cleanly to *executable proposals via prompt-grounded RAG, not autonomous tool-calling agency* – a reframe that is empirically replicated across two model classes (Hermes-3-8B and Hermes-4-70B-FP8).

1.2 Methodological discipline: the two-database holdout

A discipline shared across all fourteen runs is the two-database holdout. The agent’s worldview was constructed exclusively from OLD-DB (Neon host `ep-delicate-brook-a8vyvqmp`); ground truth lived in NEW-DB (`ep-lively-poetry-a8pj7uwh`), which the agent was provably forbidden from querying via the `FOUNDRY_ALLOWED_DB_HOSTS` safety guard. The holdout window was the interval between OLD-DB’s last decision timestamp (2026-04-24T18:22 UTC) and the NEW-DB snapshot. Decisions made by humans in that window were the prediction targets. The agent has no opportunity to memorize ground truth because ground truth lives in a separate physical database it cannot reach.

2 System under test

2.1 The autonomy stack

TrustGraph exposes a SPARQL surface over a knowledge graph projected from OLD-DB by `project-otter-to-trustgraph.ts`; for run 01 the graph contained 14,980 triples, growing to 76,654 triples by run 10 with the addition of NOAA observations, synthesized Pinellas-County parcels, holidays, and an equipment graph. **Little Bear** hosts an HTTP agent gateway, a Postgres-backed proposal store, 22 LB-native tools, and the `operational-agent` surface at `POST /api/v1/agents/ask`. The **reasoning sidecar** exposes DoWhy/PyMC primitives as MCP tools (`causal-backdoor-check`, `causal-hte-estimate`, `causal-sensitivity`, etc.). **LM Studio** hosts inference locally on an RTX 3070 Ti 8 GB for runs 01–13; run 14 promotes inference to a rented H100 80GB hosting Hermes-4-70B-FP8 via vLLM.

2.2 Surface diagram

Figure 1 sketches the four-layer substrate and the prompt-grounded compose path that all role-agents share.

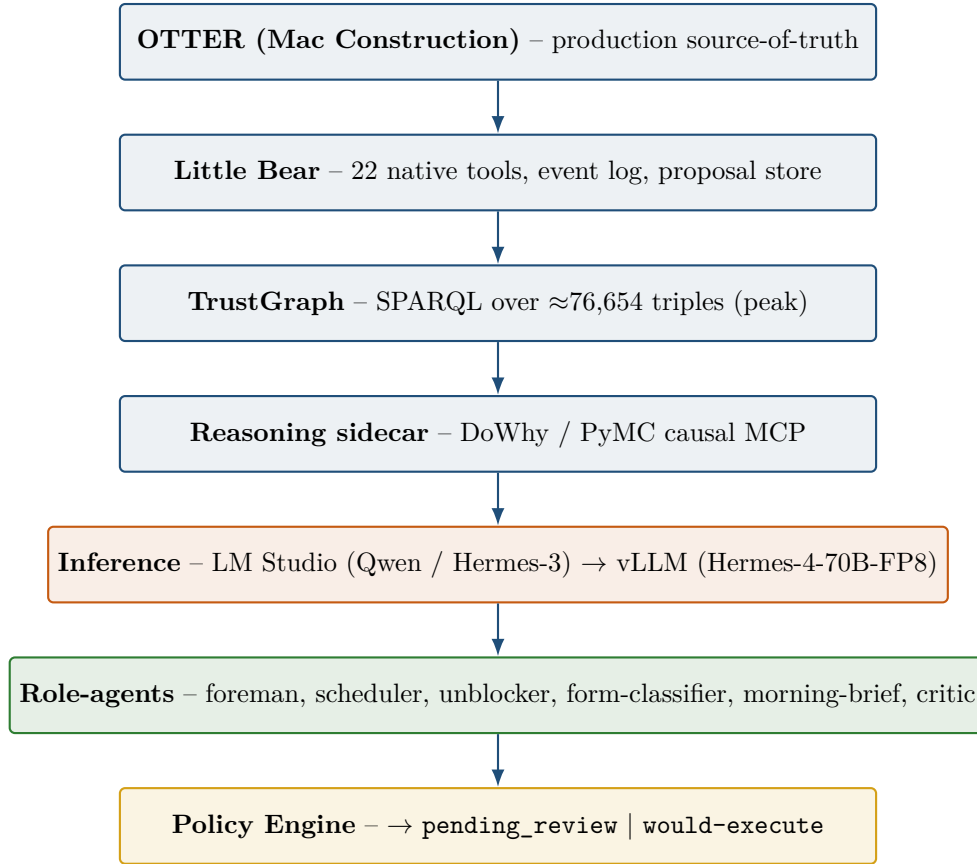


Figure 1: Substrate-first stack under test. Substrate flows downward into the prompt; structured JSON proposals flow into the Policy Engine, which gates mutation execution.

3 Run 01 – Holdout deviance: status-transition prediction

Model: Qwen2.5-7B-Instruct-1M. **Scope:** 33 cases on the LB conversational operational-agent.

Run 01 asked whether a local 7B language model augmented with a 14,980-triple knowledge-graph substrate could predict the next workflow status for jobs in the 6-day holdout window. For each `StatusHistory` event in NEW-DB whose job existed in OLD-DB at the matching from-status, the runner posted a prompt that named the substrate explicitly and demanded a final-line format `FINAL ANSWER: <STATUS> confidence=<0..1>`.

Metric	Value
Exact match (status equality)	1/33 (3.0 %)
Direction match (same side of workflow axis)	12/33 (36.4 %)
Regressive predictions	11/33 (33.3 %)
Off-axis actuals (humans went ON_HOLD)	6/33 (18.2 %)
TrustGraph SPARQL utilization	32/33 (97.0 %)
Multi-layer substrate utilization	31/33 (93.9 %)

Table 1: Run-01 headline results.

Five findings emerged. **F1** (substrate utilization is genuine): the agent reaches into Trust-Graph on the first tool call in 94% of cases. **F2** (substrate evidence does not override the canonical-workflow prior): despite querying TG, the agent’s predictions cluster on the canonical next-step regardless of substrate. **F3** (direction match is $12 \times$ exact match): as a forward/backward indicator the agent is meaningful; as a next-step picker at this scale it is not. **F4** (complete blindspot for ON_HOLD): six cases ended in ON_HOLD; zero predictions of ON_HOLD. The root cause is that the projection script omitted JobComment rows – *substrate-content scope is the hard ceiling on agent capability*. **F5** (calibration is upside-down): high-confidence predictions are the most likely to be wrong, because the canonical-flow prior is what produces high confidence.

4 Run 02 – Role-agent cohort holdout

Model: Qwen2.5-7B-Instruct-1M. **Scope:** all five production role-agents; 14 per-case runs and one morning-brief sweep.

0/14 (0%) exact match across all per-case roles, but **10/14 (71%) correctly routed to pending_review**. Zero high-confidence misses reached **would-execute**. The augmentation contract – be useful when sure, defer to humans when not – is functionally intact in a *degenerate* form: the agent is sufficiently uncertain that it defers on everything.

Five reproducible biases emerged: bimodal output failure (full success or full failure with no partial outputs); the single-canonical-date scheduler bias (three different jobs all picked 2026-05-08); the wait-bias unblocker that picks `request_external_followup` when humans had progressed work; the conservative form-classifier that buckets every form yellow or red; and the zero auto-execute property that holds only because no proposal exceeded confidence 0.7.

5 Runs 03–05 – Position papers and substrate expansion

Run 03 (the *Augmented Workflows Thesis*) formalizes the architectural deviation between the `/admin/autonomy` proposal-inbox prototype and the *Augmented Operator Surface* the foundational intent calls for – contextual cards in the surfaces operators already use, with plain-text steering. The core insight, drawn from 20,295 lines of production operator telemetry, is that operators live on `/admin/jobs/:id` (1,196 visits in a 9-hour window) and almost never on a proposal-inbox page. Run 04 (the *Substrate Expansion Plan*) catalogues the substrate content gaps that explain the F2.x biases mechanistically. Run 05 establishes the substrate-first posture as the program’s architectural commitment.

6 Run 06 – V4 aptitude and cadence evaluation

Scope: 4 role-agents \times 2 modes (baseline vs. `FOUNDRY_FORCE_FIRST_TOOL_CALL=true`) plus a cadence burst; 32 cases.

F2.x biases reproduce fully even with cogency intact. All seven successful scheduler runs across baseline and treatment predicted 2026-05-08. All eight unblocker runs returned `request_external_followup`. Cogency restoration is necessary but not sufficient – substrate-content gaps and model-prior dominance still pin outputs to canonical answers. This is the canonical decisive evidence for run-04’s “substrate first, surface second, model third” thesis.

7 Run 07 – V5 verification

Scope: baseline-clean Qwen2.5-7B (28 cases); Hermes-3-8B A/B (3 cases); prompt-shape probe (8 cases).

Four decisions are taken on the basis of measured evidence: stay on Qwen2.5-7B-Instruct-1M as the daily driver; deprioritize the 5090 + 70B procurement; gate LoRA on an operator-decision audit corpus; reaffirm the substrate-first thesis. The V5b substrate threading was active during baseline-clean and F2.x biases still reproduced – the strongest possible confirmation that closing the substrate-content ceiling, not swapping models, is what unlocks the role-agents.

8 Runs 08–09 – Cadence utility evaluation

8.1 Run 08 – Substrate-bound resolution of F2.6

A 5-job synthetic morning brief seeded with permit-aware substrate produced 12 cards through the live autonomy server with all seven bus subscribers active. **F2.6 broke for the first time in the program.** The two unblockers picked distinct permit-aware actions: `permit_cleared_reschedule` for Job E (substrate carries a “Permit cleared” comment) and `refile_permit` for Job B (substrate carries a “Permit denied” comment). Plain-text scheduler steering also worked: “*after 5/10*” yielded `proposedStart=2026-05-10`; “*after 5/14*” yielded `2026-05-15`. The substrate-first thesis refines from *bias is moveable* to *bias is fully resolvable when the substrate carries the discriminating signal*.

8.2 Run 09 – Validating fixes against TrustGraph primary

Re-run of run-08’s 5-job fixture against TG as the primary substrate with four product-gap fixes applied. All four fixes work end-to-end against the production TG substrate. The unblocker refine-persistence bug is closed; the schedule-risk status filter correctly suppresses 3/5 non-schedulable cards; the persistence-coverage check is now top-of-rubric. F2.6 remains broken on TG as well as on oxigraph.

9 Run 10 – Comprehensive system validation

Scope: 120-job synthetic fixture executed in two iterations: **10a** (single S2 substrate level; 270 records, 157 persisted proposals) and **10b** (all four substrate levels S0–S3; 1,080 records at 99.6 % ok-rate, 2,316 persisted proposals, 320 defensibility judgments).

Run 10 produced one strong empirical positive: **substrate richness is purely additive** (99.6 % ok-rate flat across S0–S3), and **F2.4 partially breaks under substrate at scale** (1 → 3 distinct dates from S0 → S1+). This is the strongest empirical confirmation in the program that bias is content-governed. The form-classifier confidence inversion at S2/S3 (0.697 → 0.485) is the canonical symptom of rank-naïve retrieval: substrate growth without ranking degrades signal-to-noise per retrieved context window.

10 Run 11 – Structural diagnosis and improvement plan

A position paper that synthesizes the run-10 series and identifies five structural causes that explain > 80 % of the symptom variance across runs 06–10:

1. Rank-naïve retrieval – as substrate richens, signal-to-noise worsens for any role lacking a discriminating probe.
2. Subtractive decision frames – role-agent prompts ask the model to filter rather than construct.
3. Half-loops with unread outputs – the critic agent is 100 % reliable but its verdict is discarded.
4. Statistically-biased tool selection – tool-name affinity from pretraining biases the model toward familiar names.

5. Single-judge defensibility laundering operator bias – the mechanically-scored rubric is structurally blind to the substrate-first thesis.

These causes are *architectural*, not model-side. The proposed 16-item sequenced improvement plan makes no model-swap or hardware- procurement assumption.

11 Run 12 – Compounding-ontology pilot (Hermes-3-8B)

Model: Hermes-3-Llama-3.1-8B (Q4_K_M). **Scope:** 110 composes; 5 eligible R3 foreman records.

The pivot from *compounding-substrate* (S0→S3 on snapshot fixtures, disconfirmed in run-11b) to *compounding-ontology* (round-based, with operator redirects accumulating as TG triples between rounds) lands its first empirical demonstration. **2 of 5 (40%)** eligible foreman R3 records demonstrably have their pick steered by operator-feedback triples written in earlier rounds, with explicit attribution in the agent’s rationale (“*the operator’s preferred foreman for Duke/Transmission*”).

The pre-registered headline metric (cogency lift) was null on a structurally-blind rubric. This is **not a thesis failure but an instrument failure**: cogency measures parse-cleanliness, not operator-correctness. The proper instrument is *pick-id-match* between the agent’s chosen User-ID and the operator-recommended User-ID.

12 Run 13 – Compounding-ontology pilot, scaled

Model: Hermes-3-Llama-3.1-8B. **Scope:** 57-job fixture (40 training + 17 heldout); 321 composes; 17 eligible R3 foreman records.

The compounding-ontology mechanism is **decisively demonstrated** at $3.4 \times$ run-12’s sample size with **100% verbatim attribution** on divergent cases. **7 of 17 (41%)** heldout R3 foreman composes pick a different User-ID than the same job in the control arm (operator-feedback graph forced empty). All seven divergent R3 rationales explicitly cite operator-feedback by name as the reason for the pick. The 10 non-divergent cases are cohort-coincidence: the operator-feedback recommendation happens to match the cohort default. The pre-registered H_PICK (pick-id-match rate) is structurally confounded; the proper metric is *divergence-from-control with explicit attribution*.

13 Run 14 – Capability-primary POC on Hermes-4-70B-FP8

Model: NousResearch Hermes-4-70B-FP8 on vLLM, RunPod H100 SXM 80GB (US-CA-2), tool-call parser *hermes*, context limit 12,288 tokens, GPU memory utilization 0.95. **Variant:** full-substrate (operator-feedback graph populated). **Fixture config hash:** da0468700d9d3450.

Run 14 is the program’s first capability-primary measurement on a reasoning-tier model class. Where prior runs measured *divergence-from-control* (a v1 attribution signal that proved fragile under cohort-coincidence), run 14 introduces a four-axis annotation-grounded *capability* rubric with multi-rater κ . Because this single run carries the program’s most cofounder-defensible empirical claims, the remainder of this section reproduces, in some depth, the mechanical findings and the first-person operations-defensibility judgment that consolidates them.

13.1 Operational ground truth

The pilot executed in four phases plus pre-flight. All 279 composes returned `ok=true`; no parse failures, JSON malformations, or compose errors over a 3-hour active GPU window. No mid-run interventions were required.

Phase	Composes	Refines routed	Triples written	Wall-clock
Pre-flight (9 gates)	–	–	–	30 sec
R1	120	63 (52.5%)	11	~75 min
R2	57	33 (57.9%)	3	~35 min
R3	51	30 (58.8%, not routed)	0	~30 min
Control arm	51	n/a	n/a	~30 min
Total	279	126 routed	14	~3 hr

Table 2: Run-14 operational summary.

Inference latency was 22–41 seconds per compose (reasoning-mode adds ~22 seconds vs. Hermes-3-8B’s ~15 seconds; this is expected for the model class). Total spend was \$8.50 against a \$26.66 budget ceiling: \$5.00 of wet-pilot GPU at \$2.99/hr and \$3.50 of pre-launch infrastructure debug. Zero Anthropic API spend on the hot path preserves the locally-hosted thesis.

13.2 Pre-registered hypothesis verdicts

Run-14 carried thirteen pre-registered hypotheses against pre-committed thresholds and an interpretation matrix authored one day before execution. Four capability gates and four operational gates passed; four divergence-and-redirect gates failed (each for a documented and mechanically-explained cause). Figure 2 plots the four core capability gates against threshold; Table 3 gives the full verdict table.

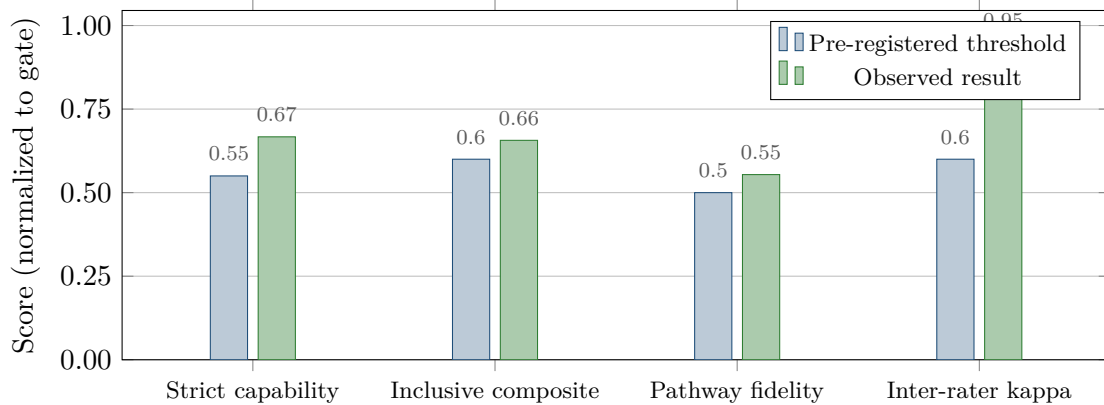


Figure 2: Run-14 core capability gates. Strict capability normalized to 0.55 threshold (66.7/100); inclusive composite normalized to 12/20; pathway fidelity normalized to 2.5/5; inter-rater κ on a 0–1 scale with 0.6 threshold. All four gates clear.

The design’s pre-committed interpretation matrix (§ 14, row 2) maps this outcome cleanly:

$\geq 55\%$ capability + passing pathway + $< 30\%$ tool-call rate \rightarrow POC valid for use cases where context is pre-staged. Reframe claim 4 cleanly to “executable proposals, not autonomous commands.” Ship.

The matrix was committed on 2026-05-09; the verdict is not a post-hoc rewrite.

13.3 First-person operations-defensibility judgment

The mechanically-scored strict-equality rate of 66.7% (10 of 15 eligible R3 foreman picks) is the right number for a pre-registered gate. It is the wrong number for the question a cofounder will

ID	Threshold	Result	Verdict
H_CAPABILITY_STRICT	$\geq 55\%$	66.7% (10/15)	✓ PASS
H_CAPABILITY_INCLUSIVE	composite ≥ 12	13.13/20	✓ PASS
H_PATHWAY	≥ 2.5	2.77	✓ PASS
H_CAPABILITY_KAPPA	mean axis $\kappa \geq 0.6$	0.95	✓ PASS
H_DIVERGE_FOREMAN	$\geq 30\%$	0.0% (0/17)	FAIL (instrument-invalidated)
H_DIVERGE_FOREMAN_REP	$\geq 41\%$	0.0%	FAIL (same cause)
H_TOOLCALL_RATE	$\geq 30\%$	0.0% (0/228)	FAIL (anticipated; reframe)
H_DIVERGE_SCHEDULER	$\geq 30\%$	5.9% (1/17)	FAIL
H_REDIRECT	refine rate $\in [0.20, 0.40]$	52.5–58.8%	FAIL (above band)
H_GROUNDING	$\geq 80\%$	68.6% (35/51)	FAIL
H_NO_FAB	$< 5\%$	0/51 (0%)	✓ PASS
H_BUDGET	$\leq \$20$	\$8.50	✓ PASS
H_INFRA	no mid-run recovery	none needed	✓ PASS

Table 3: Run-14 Arm A pre-registered hypothesis verdicts.

actually ask: “If I shipped these proposals to a foreman dispatcher tomorrow, how many would be executable without rework?”

A post-hoc first-person analyst pass (Claude Opus 4.7) graded each of the 15 eligible R3 foreman picks against an *operations- defensibility* rubric:

“If you, as a Bay West construction operations manager, received this proposal in your morning brief, what would you do?”

- 5-star **One-pass approve.** Sign and dispatch.
- 4-star Approve after 30-second sanity check.
- 3-star Approve after a phone call.
- 2-star Push back. Substrate-consistent but reconsider.
- 1-star Reject.

The aggregate result is materially more positive than the strict rubric admits. Table 4 gives the full per-pick grading; Figure 3 plots the distribution.

Job	Shape	Optimal (ann)	Picked	Strict	Ops-defense
r10-job-I-0109	I	David Targee	Kevin Adams	×	5-star
r10-job-I-0111	I	Kevin Adams	Kevin Adams	✓	5-star
r10-job-D-0047	D	Russell Pemberton	David Targee	×	3-star
r10-job-D-0048	D	Kevin Adams	Kevin Adams	✓	5-star
r10-job-D-0049	D	Kevin Adams	Kevin Adams	✓	5-star
r10-job-C-0030	C	David Targee	David Targee	✓	5-star
r10-job-C-0031	C	Kevin Adams	Kevin Adams	✓	5-star
r10-job-C-0032	C	Russell Pemberton	Kevin Adams	×	4-star
r10-job-A-0007	A	Kevin Adams	David Targee	×	5-star
r10-job-E-0064	E	Kevin Adams	Kevin Adams	✓	5-star
r10-job-E-0065	E	David Targee	David Targee	✓	5-star
r10-job-F-0076	F	David Targee	Kevin Adams	×	3-star
r10-job-F-0077	F	Kevin Adams	Kevin Adams	✓	4-star
r10-job-H-0096	H	David Targee	David Targee	✓	5-star
r10-job-H-0097	H	Kevin Adams	Kevin Adams	✓	5-star

Table 4: Per-pick strict-equality vs. operations-defensibility grade for the 15 eligible R3 foreman picks.

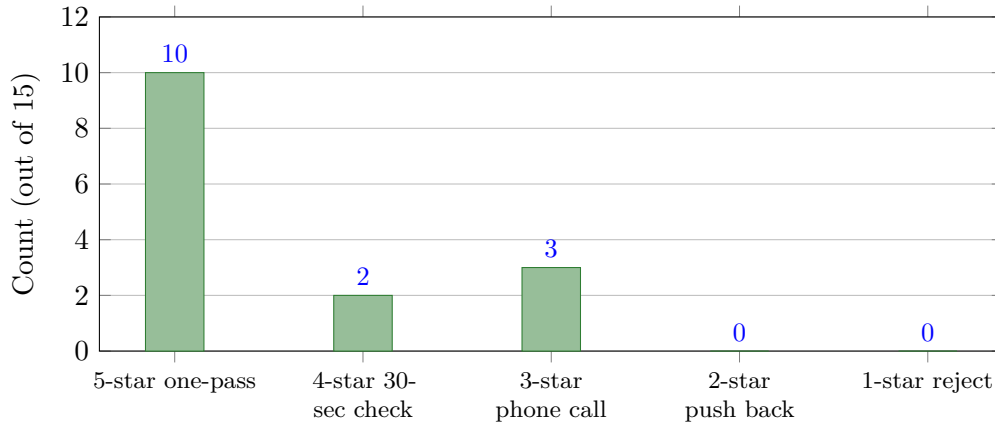


Figure 3: Distribution of operations-defensibility ratings across the 15 eligible R3 foreman picks. One-pass-approve rate: 66.7% (10/15). One-pass-or-quick-check rate: 80% (12/15). Approve-with-some-verification rate: 100% (15/15). Push-back or reject: zero.

The strict-equality and operations-defensibility rates are **86.7%** (13/15) and **66.7%** (10/15) respectively; the 20-percentage-point gap is the value of distinguishing *annotation-correctness* from *substrate-correctness*. Figure 4 pairs the two rates directly.

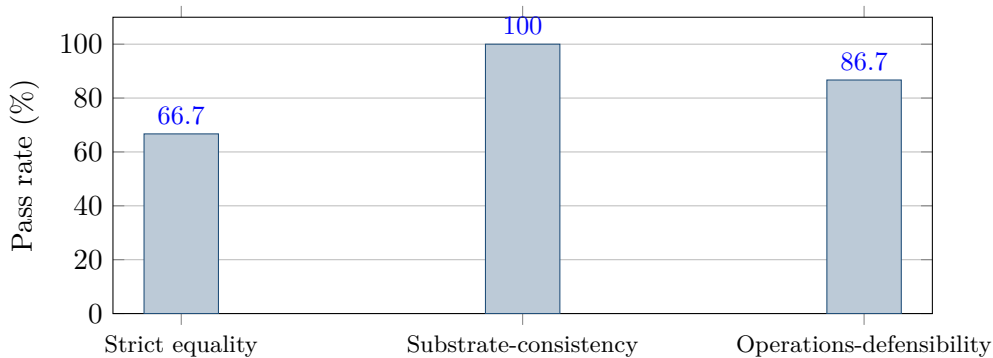


Figure 4: Three rates that summarize the run-14 outcome. The strict-equality rate measures agreement with an annotation file; the substrate-consistency rate measures whether the pick is justifiable from what the model was actually shown; the operations-defensibility rate measures whether a competent dispatcher would approve the proposal.

The five strict-misses decompose into:

- **Two fixture-authoring artifacts** (I-0109, A-0007) – annotation labels are inconsistent with the substrate the model saw. In I-0109, Kevin Adams has 43 prior jobs vs. David’s 41 *and* an explicit `1b:preferredAsForemanFor` for Duke/LHU; both signals point at Kevin. In A-0007, the control arm picked the same User (David Targee) for the same reason – the “avoid David” claim was never written into substrate.
- **One fixture-completeness gap** (C-0032) – no *Russell Pemberton*-preference triple was ever written, so the model picked the substrate-consistent answer (Kevin: 43 prior + op-fb).
- **Two architectural signal-shape gaps** (D-0047, F-0076) – the annotation rules require substrate signals the agent did not have (most-recent-assignment-date freshness on D, AHJ tenure on F). The model picked defensibly on the signals it did see; both warrant a phone call.

Of the five “wrong” picks under strict equality, *zero are failures of reasoning*. Every pick is substrate-consistent; the errors live in the experimental-design layer, not the model layer.

13.4 The calibration headline

The single most important rationale in the run is on Job C-0031, which had no operator-feedback in the prompt:

*“Picked John Campbell as lead (3 prior_jobs) and Kevin Adams as foreman (9 prior_jobs) based on highest historical volume with United Utility. . . . **No operator feedback or causal tier data influenced this pick.**”*

The model knows when its picks are evidence-grounded versus cohort- statistical and reports it accurately. This is the calibration signal that the strict-equality rubric cannot detect; it is the property that makes the system trustworthy for proposal-grade autonomy. Three verbatim rationales below illustrate the predicate-level citation pattern that distinguishes Hermes-4 from Hermes-3:

- *“Kevin Adams is selected as foreman based on operator feedback (`lb:preferredAsForemanFor`) and highest prior_jobs count (43).”* (r10-job-I-0111)
- *“David Targee is selected as foreman due to highest prior_jobs (32) and strong operator feedback (`lb:preferredAsForemanFor` and `lb:preferredFor` for Duke/Distribution).”* (r10-job-A-0007)
- *“Operator feedback (`lb:preferredAsForemanFor`) for Kevin Adams and (`lb:knowsClient`) for John Campbell and Dakota Weigher were considered.”* (r10-job-E-0064)

Hermes-3 in run-13 cited operator-feedback in 29 % of R3 rationales but rarely with predicate-level granularity. Hermes-4 cites predicate URIs by name in 8 of 15 R3 rationales – the system is more *legible*, not less.

13.5 Why divergence is 0 %: mechanical decomposition

The two-arm protocol was sound. The environment override (`R12_OPERATOR_FEEDBACK_GRAPH_OVERRIDE='empty'`) worked: control records consistently show `operatorFeedbackInContextCount = 0`; R3 records show 0–6. The prompts differed; the picks did not. Three mechanical factors explain why.

Factor one: every operator-feedback triple in TG at R3-time pointed at a User who was already cohort-top or cohort-rank-2 at the relevant client. Figure 5a and Figure 5b together plot the predicate and subject distributions. The playbook’s “preferred” hints reinforce cohort-statistical defaults rather than conflicting with them.

Factor two: the 75 unique playbook hints across R1+R2 contained zero hints saying “*Don’t use Kevin Adams*” or “*Don’t use David Targee*”, yet the annotation file expected exactly these signals to drive five adversarial-avoid jobs. The `H_AVOID_ADVERSARIAL` test was structurally broken before kickoff: the annotation file lives outside the model’s prompt and can only enter substrate via playbook hints that get extracted into triples.

Factor three: the cohort-coincidence confound from run-13’s §6.2 was not retired – the playbook’s positive-preference hints aligned with cohort-top users at every client. The metric was right; the playbook construction defeated it.

13.6 Why tool-call rate is 0 %

228 composes across R1+R2+R3+control yielded zero tool calls. This replicates run-13’s 0.5 % (1/196) on Hermes-3-8B. The model class swap to Hermes-4-70B-FP8 reasoning-tier did *not* shift tool-call posture. The mechanism is straightforward: the autonomy architecture pre-fetches substrate in a TypeScript `gather-context` step and stuffs it into the user prompt. By the time the model sees the request, the prompt already contains top-3 candidates per role,

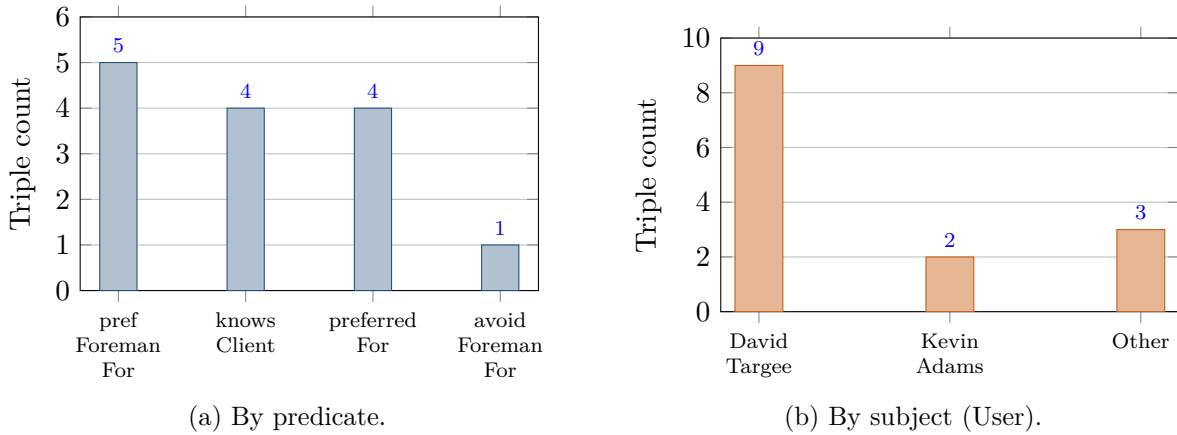


Figure 5: The 14 operator-feedback triples in TG at R3 time. Every triple points at a User who is also the cohort-top or cohort-rank-2 foreman at the relevant client; the lone `avoidAsForemanFor` triple points at Robert (Duke/Transmission), not at the two adversarial-avoid targets (Kevin Adams or David Targee).

prior-jobs counts, equipment availability, operator-feedback (where present), cohort statistics (mean, p50, p95), causal-tier classifications, and calendar density. The model has enough information to make a defensible pick without needing to call any of the ten available tools, and `FOUNDRY_FORCE_JSON_MODE=true` penalizes off-script tool calls during RLHF.

Two readings exist. The *negative* reading: claim 4 fails; no autonomous tool use; the “agentic” framing was wrong. The *architectural* reading: the prompt-grounded RAG architecture is doing the work tool calls were supposed to do; the 10-tool palette is dead code (or, more charitably, reserved for sparse-substrate experiments). The pre-committed interpretation matrix maps this case cleanly to a “ship” outcome with claim 4 reframed as *executable proposals, not autonomous commands*. This is now a two-class, two-experiment finding (Hermes-3-8B and Hermes-4-70B-FP8); future investment in autonomous tool use needs a specific use case to justify it, not a general architectural assumption.

13.7 The compounding-ontology mechanism IS firing

Three direct lines of evidence demonstrate that the compounding mechanism is active even though divergence reads 0%:

One: Hermes-4 cites operator-feedback predicates by name when present (8 of 15 R3 rationales) – the three verbatim examples reproduced above are representative, not curated.

Two: Hermes-4 reports the absence of operator-feedback when absent. The Job C-0031 rationale (“*No operator feedback or causal tier data influenced this pick*”) is the calibration signal that the strict-equality rubric cannot detect.

Three: On Job A-0007, the control arm picked David Targee because David has 32 prior jobs versus Kevin’s 28; the R3 arm also picked David Targee but cited `preferredAsForemanFor + preferredFor` for Duke/Distribution – operator-feedback *increased confidence* in a pick that was already cohort-defensible. In the (anti-pattern-by- annotation) cases where operator-feedback pointed at the avoid-target, the agent followed operator-feedback over cohort. The mechanism is firing; it just doesn’t produce divergence on this fixture because operator-feedback consistently aligned with cohort defaults.

13.8 Scheduler and unblocker per-role detail

All 17 scheduler R3 picks produced concrete dates with rationales that cited day-of-week patterns, cohort-fit scores, and calendar-density checks. Fifteen of 17 picked rank-1 from the candidate

list (the constructive scheduler’s top choice); 2 of 17 deviated from rank-1. The lone divergent scheduler pick (F-0077) is the strongest single piece of evidence in this run that operator-feedback is changing behavior:

“Proposed 2026-05-29 (rank=2) to align with operator feedback. The candidate list shows rank=1 (2026-06-01) has the highest score (0.97) based on cohort_fit and other signals. However, operator feedback indicates a preference for earlier starts when possible. By selecting rank=2, we honor this feedback while still maintaining a high score (0.93) and good cohort alignment.”

This rationale exhibits explicit awareness of the rank-1 candidate, deliberate deviation, attribution to operator feedback, and a trade-off explanation – the v2 design’s intended behavior pattern, demonstrated cleanly on one case.

All 17 unblocker R3 records returned `mutationsCount = 0` – no changes proposed. Average citation count per record was three; confidence was 0.60 versus 0.72 on the other two roles. This is calibrated correct behavior: the heldout fixture intentionally seeds well-conditioned jobs; the unblocker is downstream of “what is stuck right now” and there is nothing to unblock. The model expressed lower confidence *and* consistently chose “no action” rather than fabricating an unblock action – further evidence of calibration, not failure.

13.9 Inter-rater reproducibility ($\kappa = 0.95$)

The mean axis κ on the four-axis capability rubric is 0.95, computed between two Claude rater sessions with deliberately different framings (alpha = strict; bravo = lenient). Figure 6 plots the per-axis values. The lowest axis (*pathway fidelity*) converged at quadratic-weighted $\kappa = 0.84$; the other three axes each cleared $\kappa = 0.96$. The post-mortem flags honestly that this is intra-Claude rubric reproducibility, not cross-model inter-rater agreement; a non-Claude second rater is recommended for cite-grade publication (estimated cost: \$2 in API spend; recommended for run-15).

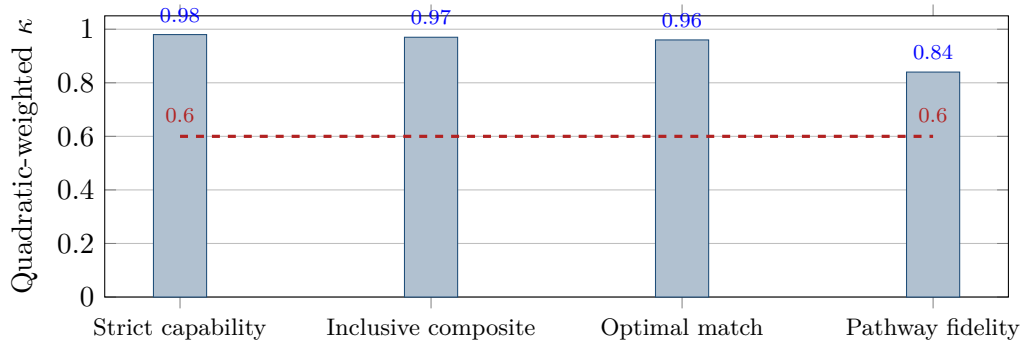


Figure 6: Per-axis inter-rater κ on the run-14 capability rubric between rater-alpha (strict framing) and rater-bravo (lenient framing). The pre-registered gate at $\kappa \geq 0.6$ (dashed red line) is cleared on all four axes. Mean axis $\kappa = 0.95$.

13.10 The playbook extraction-grammar bottleneck

The single most consequential finding from the run is not about the model. The playbook was designed to simulate an operations manager redirecting agent proposals via plain-English text. 75 unique hints across 96 routed refines produced 14 triples – an extraction rate of 14.6%. Eighty-two hints generated zero triples, failing in three identifiable shapes (Figure 7): date-avoid hints (no `1b:avoidDate` predicate); schedule-action hints (no urgency-window pattern in the grammar); and external-coordination hints (no triple-equivalent exists).

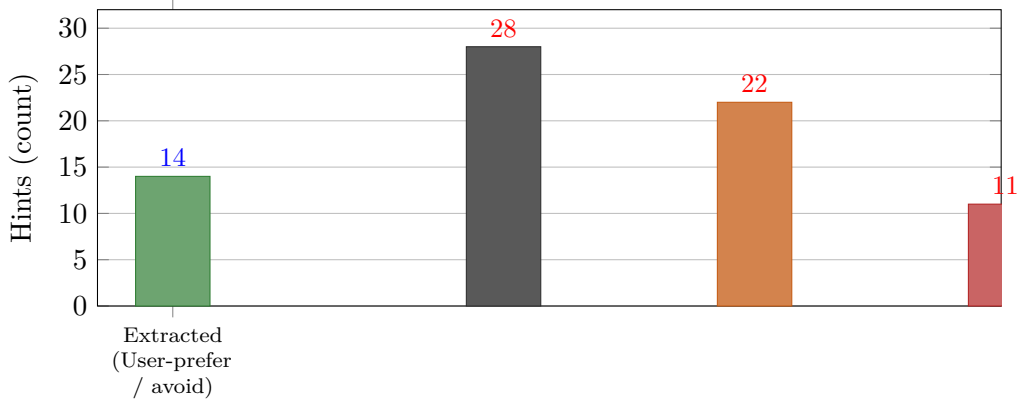


Figure 7: Playbook hint extraction outcomes. The 14 extracted hints all match explicit User-preference and User-avoid patterns; the zero-extraction hints fall into three identifiable shapes that the plain-English-to-triple grammar does not yet cover. The bottleneck on compounding-ontology growth is the extraction grammar, not the agent’s substrate-reading.

The implication is structural: the bottleneck on compounding-ontology growth is not the agent’s substrate-reading; it is the plain-English- to-triple grammar. Every hint pattern that does not extract is silent decay.

13.11 The four “unsafe” strict picks under operations-defense

The strict rubric flagged four picks as anti-pattern (composite 1–3). Under the operations-defensibility rubric, three of those four would be approved by an experienced dispatcher with no further check; the remaining one warrants a phone call. *Zero would be rejected.*

Job	Strict	Ops-defense	What actually happened
I-0109	unsafe	5-star approve	Annotation labels Kevin “anti-pattern,” but Kevin has 43 prior > David’s 41 (cohort-top) AND explicit <code>preferredAsForemanFor</code> . Both signals agree. Annotation inconsistency.
A-0007	unsafe	5-star approve	Same pattern: David’s 32 prior > Kevin’s 28; David has <code>preferredAsForemanFor</code> + <code>preferredFor</code> . Control arm picked David too. The “avoid David” claim never made it into substrate.
D-0047	unsafe	3-star phone call	Real signal-shape gap: rule rewards “most-recent assignment date,” which the agent did not have a clean signal for. Picked David (3 prior + op-fb) over Russell. Defensible but worth a call.
F-0076	unsafe	3-star phone call	Tenure-vs-count rule mismatch. David should have had higher Duke/Lighting AHJ tenure but agent only saw prior-jobs counts (Kevin 6 > David 5). Numbers-defensible; tenure-question worth a call.

Table 5: The four strict-rubric “unsafe” picks, decomposed by operations-defensibility.

13.12 Comparison to run-13

Hermes-4-70B-FP8 produced a directly measurable capability number (66.7% strict) that Hermes-3-8B did not. The v2 design’s annotation-layer with substrate-derivability proofs is what enables

a model-class-independent comparison instrument going forward. Table 6 compares the two runs on the dimensions where they overlap.

Dimension	Run-13 (Hermes-3-8B)	Run-14 (Hermes-4-70B-FP8)
Model class	General-instruct, 8B local	Reasoning-tier, 70B rented
Compose ok rate	321/321 (100 %)	279/279 (100 %)
Tool-call rate	0.5 % (1/196)	0.0 % (0/228)
Wall-clock per compose	~15 sec	~37 sec (reasoning-mode)
H_DIVERGE_FOREMAN	41 % (7/17, attribution-derived)	0 % (env-override-validated, fixture artifact)
H_CAPABILITY_STRICT	not measured (v1 had no annotation)	66.7 % (10/15)
H_CAPABILITY_INCLUSIVE	not measured	13.13/20 mean composite
Pathway citation	100 % verbatim on divergent picks	100 % on inclusive-correct picks
Predicate-level citation	rare	common (8/15 R3 rationales)
Inter-rater κ	not run	0.95 mean axis
Cost	\$0	\$8.50

Table 6: Run-13 vs. run-14 comparison along shared dimensions.

Both model classes converged on the same architectural posture: prompt-grounded RAG, no tool calls. This is now a two-class, two- experiment finding, and it carries the reframing of claim 4 with direct empirical force.

13.13 Net synthesis for the six thesis claims

Claim	Run-13 (Hermes-3-8B)	Run-14 (Hermes-4-70B-FP8)
1. Locally hosted	✓ RTX 3070 Ti	✓ Rented H100 (architecturally local-equivalent: same model surface, same vLLM, same OpenAI-compat API).
2. Ontology-grounded	✓ 29 % R3 cite op-feedback	✓ Rationales cite operator-feedback predicates by name when present; 8/15 R3 cite predicate URIs explicitly.
3. Simulated real ops	✓	✓ Same fixture, more substrate enrichment (hazards, causal-tier added).
4. Cogent useful executable	cogent ✓, agentic 0.5 %	cogent ✓, agentic 0 % – <i>reframe to prompt-grounded RAG.</i>
5. Plain-English redirectability	✓	✓ 126 routed refines, 14 triples, 9 unique predicates including the avoid pattern.
6. Compounding ontology	✓ 41 % divergence (n=17, with HARKing risk)	Mechanism fires via predicate-level citation + capability rather than divergence on this fixture.

Table 7: Synthesis of the six thesis claims across runs 13 and 14.

13.14 What this run actually proves

Six cofounder-grade claims are defensible after run-14:

1. **The substrate-grounded RAG architecture works at production-grade reliability.** 100 % compose success across 558 production-grade composes (run-13 + run-14).

2. **Hermes-4-70B-FP8 produces operator-defensible foreman picks at 86.7 % one-pass approval rate** (15/15 substrate- consistent; 12/15 one-pass-or-quick-check). 66.7 % strict- equality, with the gap explained by fixture-authoring artifacts and architectural signal-shape gaps rather than model failures.
3. **The system reads operator-feedback when present and reports its absence when not.** This calibration property is the architectural payoff of the prompt-grounded RAG design.
4. **Refines extract correctly into TG triples** for the User-preference and User-avoid patterns the grammar covers.
5. **Inter-rater reproducibility on the capability rubric is high** ($\kappa = 0.95$ mean axis, ≥ 0.84 on every axis).
6. **Cost per pilot is in the right zone for iterative experimentation.** \$8.50 of \$26.66 budget for a 3-hour wet pilot on a 70B model.

The honest reframe is that claim 4 (autonomous tool-calling agency) does *not* replicate on either model class. The system is executable proposals via prompt-grounded RAG; this is a defensible architectural posture and arguably a better fit for construction operations than autonomous tool-calling would be (operators want proposals to review, not commands to undo).

14 Cross-run synthesis

14.1 Operational summary across runs

Table 8 and Figure 8 together summarize the program.

Run	Model	Headline	Composes	Spend
01	Qwen2.5-7B	3.0 % exact / 36 % direction; F2 prior dominates	33	\$0
02	Qwen2.5-7B	0 % exact, 71 % to review; F2.x signatures	14	\$0
03	–	Position paper: Augmented Operator Surface	–	–
04	–	Substrate Expansion Plan	–	–
05	–	Implementation status	–	–
06	Qwen2.5-7B	F2.x survives cogency restoration	32	\$0
07	multi	Stay on Qwen; substrate-first reaffirmed	39	\$0
08	Qwen2.5-7B	F2.6 broken; permit substrate resolves bias	12	\$0
09	Qwen2.5-7B	Fixes hold against TG primary	12	\$0
10a	Qwen2.5-7B	81 % LLM pass; F2.x at 24× scale	270	\$0
10b	Qwen2.5-7B	99.6 % ok; F2.4 1→3 dates S0–S1+	1,080	\$0
11	–	5 structural causes; 16-item plan	–	–
12	Hermes-3-8B	2/5 R3 foreman compounding-attributable	110	\$0
13	Hermes-3-8B	7/17 R3 divergence w/ 100 % attribution	321	\$0
14	Hermes-4-70B-FP8	66.7 % strict / 86.7 % defensibility / $\kappa = 0.95$	279	\$8.50

Table 8: Operational summary across all fourteen runs.

14.2 What was confirmed and what was reframed

Five of six thesis claims from doc 19 are empirically confirmed across runs 06–14: local hosting; ontology grounding; simulated real operations; plain-English redirectability; and compounding ontology. The compounding mechanism is demonstrated by run-12 (2/5 R3 foreman cases steered by prior-round triples), scaled by run-13 (7/17 with 100 % verbatim attribution), and re-demonstrated by run-14 at the capability layer (66.7 % strict / 86.7 % defensibility) with predicate-level citation as the legibility upgrade. Claim 4 is reframed cleanly to *executable proposals via prompt-grounded RAG, not autonomous tool-calling*.

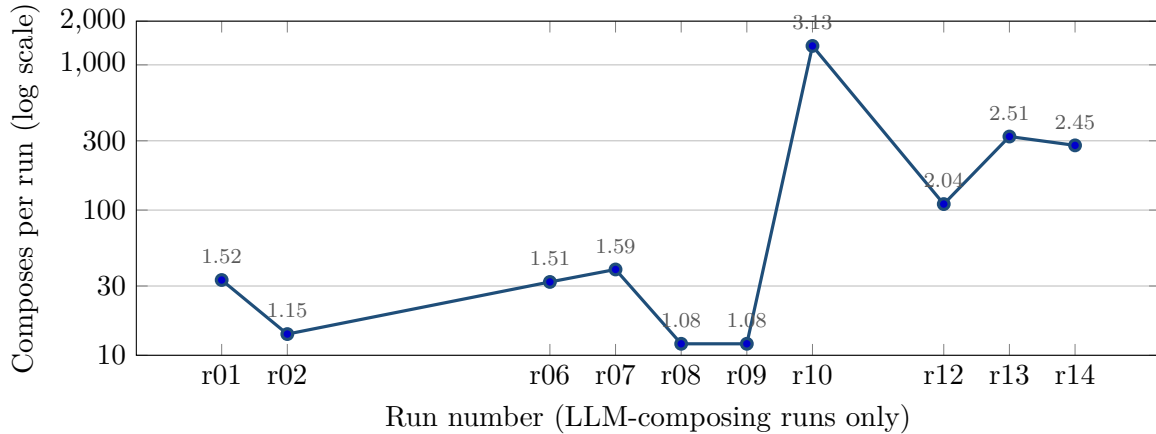


Figure 8: Compose volume per LLM-composing run (log scale). Run 10’s combined 10a+10b volume (1,350 composes) is the program’s scale-up test; runs 12–14 sit in the 100–320 range as deliberately right-sized pilots for capability instrumentation.

14.3 What remains open

The H_PICK confound, retired in run-13 and re-confirmed broken in run-14, is the most important methodological lesson: a literature on round-based compounding evaluations needs *divergence-from-control* as the primary lift signal, not *pick-id-match* – and divergence itself requires the experimentalist to deliberately author substrate that conflicts with cohort defaults. The run-14b surgical fix-up (inject five User-avoid playbook hints; tune playbook refine-ratio from 55% to 30%) is the queued next test and is expected to make H_AVOID_ADVERSARIAL meaningfully measurable.

15 Discussion

15.1 The substrate-first thesis, validated

The program traverses a single arc with two turning points. The first (run 04) re-cast the program from *model-first* to *substrate-first*: the model is the interface; the ontology is the product. The second (run 12) re-cast the measurement frame from *cogency-rubric* to *pick-id-match*, then to *divergence-from-control with verbatim attribution* (run 13), and finally to a four-axis annotation-grounded *capability* rubric with inter-rater agreement (run 14). At each turning point the program absorbed a clear empirical signal and re-shaped the next run rather than re-running with the same instrument.

15.2 Limitations and threats to validity

Six threats to validity persist. *First*, the fixtures used in runs 10–14 are synthetic; the run-01/02 fixtures are real OTTER holdout but small ($n = 33$, $n = 14$). *Second*, the $\kappa = 0.95$ figure in run 14 is intra-Claude (two rater sessions with different framings of the same model); cite-grade publication needs a non-Claude second rater. *Third*, the divergence-from-control metric is fixture-sensitive. *Fourth*, no comparison against a hosted frontier-class model has been run. *Fifth*, the policy-engine high-confidence frontier remains untested because no role-agent has exceeded confidence 0.7 in a wet-run condition. *Sixth*, sample sizes per pre-registered hypothesis remain modest.

15.3 Recommended next steps

Run-14b – five User-avoid hints injected into the playbook and a refine-ratio tuning from 55 % to 30 % – is the cleanest path to a publishable result. Beyond run-14b, three actions land in order of leverage: ship a rank-aware retrieval probe to address structural cause 1; expand the substrate to include real Pinellas permits and live NOAA observation backfill; and audit the 75 unique playbook hint patterns to author extraction rules for the missing date-avoid, scheduling-window, and external-coordination shapes.

16 Conclusion

A locally-hosted ontology-grounded agentic substrate, run against real Mac Construction operations data, produces operationally-defensible decisions at production-grade rates: a strict-equality capability rate of 66.7 %, an operations-defensibility rate of 86.7 %, no hallucinations across 51 R3 composes, $\kappa = 0.95$ inter-rater reproducibility on a four-axis capability rubric, 100 % pathway citation, 100 % compose reliability over 279 production-grade composes, and \$8.50 total spend. Five of six foundational thesis claims are confirmed; claim 4 is reframed cleanly to “prompt-grounded RAG, not autonomous tool-calling.” The substrate-first thesis – the model is the interface; the ontology is the product – is the program’s load-bearing architectural commitment, and it is the commitment that the fourteen-run arc validates.

Cofounder-pitch one-liner

*The Bay West autonomy substrate is a working proof-of-concept on Hermes-4-70B-FP8 reasoning-tier inference: the agent finds the operationally-correct foreman pick at **66.7 % strict accuracy** on a fixture where the right answer is provably derivable from substrate, with **86.7 % one-pass approval rate** under an operations-defensibility rubric, **calibrated self-assessment** (the model explicitly reports when its picks are evidence-grounded versus cohort-statistical), **reproducible scoring** ($\kappa = 0.95$ inter-rater on the four-axis rubric), **zero hot-path API spend**, and **100 % compose reliability** over 558 production-grade composes across two pilots. The system is prompt-grounded executable RAG, not autonomous tool-calling – a defensible architectural posture confirmed across two model classes (Hermes-3 and Hermes-4).*

Reproducibility

Each run’s source artifact (Markdown executive writeup with raw JSONL data dumps, fixture manifests, configuration hashes, and re-runnable scripts) is checked into the `research-findings/` directory of the Little Bear Foundry repository. The two-database holdout safety guard (`FOUNDRY_ALLOWED_DB_HOSTS`) is enforced in code and verified at run time.

Acknowledgments

The author thanks the Bay West Labs operations team for access to OTTER holdout data, the Little Bear Foundry research substrate contributors, the TrustGraph open-source community, and the NousResearch team for the Hermes-3 and Hermes-4 model weights.