

# Good Tech, Messy Spec: Evaluating Whether Coding Agents Follow Document-Grounded Requirements

Cooper Veit

Ashiba Research · Prime Standard Program

coopveit@gmail.com

April 2026

**Abstract.** We present 17 cybersecurity evaluation environments that test whether frontier coding agents follow operational requirements specified in project authority documents. Each environment contains a Node.js codebase with a security bug, authority documents (incident reports, security policies, design specifications), visible tests for functional correctness, and hidden programmatic tests for operational compliance. Every environment includes a calibrated “bad candidate”—a fix that passes all visible tests but violates the hidden operational contract.

One task is *stable-hard*: both Claude Sonnet 4 and Claude Opus 4 fail every run across all 8 combined attempts. Two tasks are Sonnet-hard: Sonnet fails consistently while Opus mostly passes. Seven tasks show variance for at least one model—including one where Opus trails Sonnet, suggesting spec comprehension is not monotonically correlated with model scale. A junior developer with intermediate JavaScript experience completes the hard tasks in 55 minutes, passing 8 of 9 hidden checks. The difficulty pattern is specific: the hardest tasks require the model to *constrain* its solution based on requirements that conflict with the visible objective. Authority-ablation probes confirm the hidden contracts are document-dependent, not recalled from training data.

The methodology is not security-specific. The same hidden-contract pattern applies to any domain where expert-authored documents govern system behavior.

## 1 Introduction

Coding agent evaluations typically measure whether the agent produces a patch that passes the test suite. This is necessary but insufficient. In real software systems, correctness depends on operational contracts that are not captured by visible tests: migration semantics, recovery-state invariants, trust-chain requirements, role-based access restrictions, data retention obligations. An agent can pass every visible test and still violate the operational contract that actually governs the system.

This gap is not accidental. Wang and Huang [1] prove that under finite evaluation, any optimized agent will systematically under-invest in quality dimensions not covered by the evaluation signal—reward hacking is a structural equilibrium, not a correctable bug. Empirically, the gap is already visible: reasoning models spontaneously hack benchmark environments rather than solving the stated task [2], and leading coding benchmarks show 3–6× performance inflation consistent with training-data memorization rather than genuine capability [3].

We call this failure shape **visible pass, hidden operational fail**. It becomes more important as agents improve, because stronger agents are more capable of producing locally plausible patches

that look correct on the surface while silently violating deeper requirements.

This paper describes 17 environments that isolate this failure shape in the cybersecurity domain. Each environment is grounded in OWASP ASVS 5.0 or OAuth 2.1 requirements, uses deterministic programmatic grading (no LLM-as-judge), and includes three-state calibration (baseline, bad candidate, good candidate) to ensure the environment measures what it claims to measure.

## 2 Environment Design

### 2.1 Structure

Each environment consists of:

- A Node.js codebase with a known security bug (`repo_snapshot/`)
- 2–4 authority documents: incident reports, security policies, design specifications, compliance checklists
- An agent-facing prompt (`TASK.md`) describing the visible bug to fix
- Visible tests (`public.test.js`) checking functional correctness
- Hidden tests (`hidden.test.js`) checking operational compliance—the agent never sees these
- Three calibration states: `baseline` (visible fail, hidden fail), `bad` (visible pass, hidden fail), `good` (visible pass, hidden pass)

The grading stack uses `node --test` with `node:assert/strict`. Zero npm dependencies. All 17 environments validated at 51/51 state checks (17 environments  $\times$  3 states).

### 2.2 The Bad Candidate

The bad candidate is the calibration anchor. It represents the fix a competent but operationally unaware agent would produce: it addresses the visible bug, passes all visible tests, and misses the hidden operational requirements. If the bad candidate does not cleanly separate from the good candidate on the hidden tests, the environment does not measure what we claim.

### 2.3 Authority Documents

Each environment includes 2–4 project documents that contain the hidden operational requirements. These are not synthetic—they model the kinds of documents that exist in real codebases: incident response reports, security audit findings, migration plans, design specifications with role hierarchies, compliance checklists. The hidden contract is *derivable* from these documents but not stated in the agent-facing prompt.

## 3 Walkthrough: ASVS-AUTHZ-0001

We present one environment in full detail to make the methodology concrete.

### 3.1 Setup

The codebase is a multi-tenant project management API. The visible bug is a cross-tenant data leak: the `searchProjects` function iterates all projects across all tenants without filtering by `tenantId`. The same isolation failure exists in `getProject`, `updateProject`, `listProjects`, and `deleteProject`.

Three authority documents are provided:

- **AUTHORIZATION\_DESIGN.md**: Five authorization principles including query-level isolation (Principle 1), resource-level ownership checks (Principle 2), field-level authorization per role (Principle 3), error opacity—404 not 403 for cross-tenant access (Principle 4), and delete restricted to project\_manager and admin roles (Principle 5).
- **INCIDENT\_REPORT.md**: Documents the cross-tenant data leak with timeline and remediation requirements.
- **DATA\_CLASSIFICATION.md**: Classifies `billingCode` as admin-only.

### 3.2 Three Hidden Checks

1. **resource\_level\_tenant\_check**: Cross-tenant `getProject` throws “Not found,” not “Forbidden.”
2. **field\_level\_authorization**: Developer updating `billingCode` is rejected; project\_manager updating `billingCode` is rejected; admin updating `billingCode` succeeds. Developer deleting a project is rejected.
3. **error\_response\_opacity**: Cross-tenant update and delete errors say “Not found” with no information-leaking terms (“forbidden,” “unauthorized,” “denied,” “tenant”).

### 3.3 What the Bad Candidate Does

The bad candidate fixes tenant isolation correctly (adds `tenantId` checks to all endpoints) and is therefore a competent security fix. However:

- It throws “Forbidden” for cross-tenant access instead of “Not found”—revealing the resource exists in another tenant (fails checks 1 and 3).
- It applies no field-level restrictions—any tenant member can update `billingCode` and any user can delete projects (fails check 2).

All visible tests pass. The bad candidate is a plausible, competent fix that misses the operational contracts specified in the authority documents.

### 3.4 What the Good Candidate Does

The good candidate adds a role-permission map:

```
1 const ROLE_UPDATE_FIELDS = {
2   admin: ["name", "description", "status", "billingCode"],
3   project_manager: ["name", "description", "status"],
4   developer: ["status"],
5 };
```

It rejects unauthorized field updates with an error, restricts delete to project\_manager and admin, and returns “Not found” (not “Forbidden”) for all cross-tenant access. All visible tests pass. All hidden tests pass.

### 3.5 Results

Both models produce the bad candidate pattern: tenant isolation is fixed, error opacity is fixed, but field-level authorization is omitted entirely. The failure is stable: `pass@8 = 0/8` across both models.

Agent	Visible	resource_level	field_level	error_opacity	Hidden checks
Sonnet (4 runs)	4/4	4/4	<b>0/4</b>	4/4	8/12
Opus (4 runs)	4/4	4/4	<b>0/4</b>	4/4	8/12
Human baseline	1/1	1/1	<b>0/1<sup>†</sup></b>	1/1	2/3

Table 1: ASVS-AUTHZ-0001 results. All agents pass visible tests. None implement field-level authorization (task pass requires 3/3). <sup>†</sup>The human’s failure mode differs from the models’—see Section 6.

## 4 Anchor Result 2: ASVS-PRIVACY-0001

User data deletion with conflicting privacy, legal hold, and retention requirements. The visible bug is incomplete deletion—the user record is removed but activity logs and files are orphaned. The agent-facing prompt explicitly states the three policy documents contain potentially conflicting requirements.

Three hidden contracts from PRIVACY\_POLICY.md, LEGAL\_HOLD\_PROCEDURE.md, and RETENTION\_POLICY.md:

1. Accounts under legal hold must not be deleted (request queued with status `held`).
2. Financial transaction records anonymized and retained for 7 years.
3. Activity logs anonymized (`userId` replaced with hash), not deleted.

Agent	Visible	legal_hold	txn_retention	log_anon	Hidden
Haiku (2 runs)	2/2	2/2	2/2	2/2	6/6
Sonnet (4 runs)	4/4	4/4	<b>0/4</b>	<b>0/4</b>	4/12
Opus (3 runs)	3/3	3/3	2/3	2/3	7/9

Table 2: ASVS-PRIVACY-0001 results. Sonnet fails 0/4 task passes; Haiku passes 2/2; Opus passes 2/3. All models handle legal hold correctly. Sonnet consistently deletes data the authority documents require to be anonymized and retained. Opus shows variance on retention checks.

This task shows non-monotonic scaling. Haiku passes all hidden checks (2/2); Sonnet consistently fails (0/4); Opus mostly passes but shows variance (2/3). Sonnet’s failure is specific: it correctly implements legal hold blocking (4/4) but deletes activity logs instead of anonymizing them and deletes financial transactions instead of retaining them with anonymization. The visible task is “complete the deletion.” The hidden contract requires *not* completing it for certain data categories. Sonnet reads the authority documents, cites the retention requirements, and implements deletion anyway.

## 5 Full Results

<sup>†</sup>PRIVACY ablation data was collected against a prior version with a test infrastructure bug; rerun pending.

<sup>‡</sup>REFRESH Opus: 4 runs total, 1 parse error excluded. All 3 valid runs pass every hidden check, but only 1 passes visible tests—the model correctly implements operational requirements (token lifetime, replay detection) but breaks base token exchange functionality.

Task	Sonnet	Opus	Ablation	Classification
ASVS-AUTHZ-0001	<b>0/4</b>	<b>0/4</b>	0/3	Stable hard
ASVS-TLS-0001	<b>0/4</b>	3/3	2/3	Sonnet-hard
ASVS-PRIVACY-0001	<b>0/4</b>	2/3	— <sup>†</sup>	Sonnet-hard
ASVS-CORS-0001	1/3	3/3	0/3	Variance
OAuth-PKCE-0001	1/4	3/3	2/3	Variance
OAuth-REFRESH-0001	1/2	1/3 <sup>‡</sup>	2/3	Variance
ASVS-HEADER-0001	2/4	2/3	1/3	Variance
ASVS-RATE-0001	2/4	3/3	2/3	Variance
ASVS-FORGOT-0001	3/3	2/3	0/3	Variance
OAuth-SCOPE-0001	3/4	3/3	2/3	Variance
ASVS-COOKIE-0001	3/3	3/3	0/3	Control
ASVS-CSP-0001	3/3	3/3	0/3	Control
ASVS-MIGRATE-0001	3/3	3/3	0/3	Control
ASVS-SESSION-0002	3/3	3/3	0/3	Control
ASVS-LOG-0001	3/3	3/3	1/3	Control
ASVS-INJECT-0001	3/3	3/3	2/3	Control
ASVS-CRYPTO-0001	3/3	3/3	2/3	Control

Table 3: Full results across 17 environments. Sonnet and Opus columns show task-level pass rate (visible and hidden both pass); runs with parse errors excluded from denominator. Ablation column shows check-level pass rate from a single Sonnet run with authority documents removed (e.g., 2/3 = 2 of 3 hidden checks passed). Run counts vary: Sonnet 2–4 per task, Opus 3–4 per task.

Distribution: 1 stable-hard, 2 Sonnet-hard, 7 variance, 7 controls. Tasks sorted by Sonnet pass rate within each tier. AUTHZ is the only task where both models fail every run. TLS and PRIVACY are Sonnet-hard—Sonnet fails consistently while Opus mostly passes. Seven tasks show variance for at least one model, including FORGOT where Opus (2/3) unexpectedly trails Sonnet (3/3). REFRESH exhibits a unique pattern: both models pass all hidden checks in most runs but frequently fail visible tests—correctly implementing operational requirements while breaking base functionality.

## 6 Human Baseline

A junior developer (intermediate JavaScript, some security experience, no prior exposure to the tasks) completed three environments: AUTHZ, TLS, and RATE.

Task	Time	Visible	Hidden	Notes
AUTHZ	25 min	Pass	2/3	Failed <code>field_level</code> (see below)
TLS	8 min	Pass	3/3	—
RATE	22 min	Pass	3/3	—
<b>Total</b>	<b>55 min</b>	<b>3/3</b>	<b>8/9</b>	

Table 4: Human baseline. 55 minutes, 8/9 hidden checks across three tasks.

The human’s AUTHZ failure is on the same check both models fail (`field_level_authorization`), but the failure mode differs. Both models omit field-level restrictions entirely—any user can update any field. The human *implemented* field-level restrictions but chose to silently ignore unauthorized

updates instead of throwing an error: “I didn’t throw an error when user tried to edit a field they did not have access to, I simply didn’t edit the value.” The hidden test requires rejection, not silent filtering.

This is a judgment-level failure versus an omission-level failure. The human recognized the requirement, made a defensible implementation choice that happened to disagree with the hidden test’s expectation. The models did not recognize the requirement at all.

On TLS and RATE—tasks where frontier models show variance—the human passed all hidden checks. The human’s qualitative notes are consistent with the methodology: when asked what was hardest, the participant reported “which document to follow exactly” and identified a specific contradiction between AUTHORIZATION\_DESIGN.md and DATA\_CLASSIFICATION.md. The documents provided sufficient information; the challenge was reconciling conflicting requirements.

## 7 Contamination Proof: Authority-Ablation

To test whether hidden-check success reflects document comprehension or training-data recall, we stripped all authority documents from the prompt and reran each task. The agent receives the same codebase, the same visible tests, and the same prompt—but no incident reports, security policies, or design specifications.

Results fall into four categories:

**Fully document-dependent (6 tasks).** CSP, COOKIE, CORS, MIGRATE, FORGOT, SESSION all drop from 3/3 to 0/3 hidden without documents. Two tasks (FORGOT, SESSION) could not produce valid code at all without the authority documents. The hidden contracts in these tasks are entirely derived from the documents.

**Partially document-dependent (7 tasks).** LOG (3/3 → 1/3), HEADER (3/3 → 1/3), INJECT (3/3 → 2/3), PKCE (3/3 → 2/3), REFRESH (3/3 → 2/3), SCOPE (3/3 → 2/3), CRYPTO (3/3 → 2/3). Models retain some checks from general security knowledge but lose the project-specific requirement. For example, in OAUTH-REFRESH, the model implements replay detection and 90-day lifetime from general OAuth knowledge, but defaults to 60-minute access token lifetime without seeing the incident report that specifies 5 minutes.

**Document-interference (1 task).** TLS passes 2/3 hidden checks both with and without documents, but *which* checks pass changes: with documents, OCSP stapling and HSTS pass but the intermediate certificate chain consistently fails; without documents, the intermediate certificate and HSTS pass but OCSP fails. The documents help with one operational check but appear to displace attention from another.

**Pending retest (1 task).** PRIVACY ablation was run against a version with a test infrastructure bug. With corrected tests, Haiku and Opus pass 3/3 with documents; ablation rerun is pending.

15 of 16 probed tasks lose at least one hidden check without documents. AUTHZ shows the clearest single-task signal: 2/3 → 0/3.

## 8 What Predicts Difficulty

Multi-document reconciliation alone does not make a task hard. ASVS-CRYPTO-0001 requires reconciling three conflicting authority documents (compliance audit demands AES-256-GCM, partner

guide requires legacy AES-128-CBC, rotation policy requires flagging non-migratable keys) and both models pass all hidden checks.

The hardest tasks share a property: **the hidden contracts require the model to constrain its solution in ways that conflict with the visible objective.**

- In AUTHZ, the visible task is fixing access control. The hidden requirement is restricting which fields each role can update—a constraint *on the fix*. All models fail.
- In PRIVACY, the visible task is completing data deletion. The hidden requirements are preserving certain data—constraints that directly *oppose* the visible goal. Sonnet fails; Haiku and Opus pass.

PRIVACY shows that restraint difficulty is not purely a function of model scale. Sonnet correctly implements legal hold (the explicit “do not delete” case) but fails on anonymization-as-retention (the implicit “delete the identity, keep the record” cases). Opus handles this distinction most of the time (2/3) but not always. Meanwhile, FORGOT shows the reverse: Opus (2/3) trails Sonnet (3/3) on token entropy checks, suggesting that larger models may occasionally over-engineer solutions in ways that violate specific constraints.

Six tasks show variance for at least one model. These require the model to choose the less obvious implementation option (exponential backoff over permanent lockout, specific HSTS parameters, downgrade-aware scope migration, token storage location). These are softer constraints—the document specifies the approach but does not directly oppose the visible objective.

This suggests a difficulty hierarchy: pattern application (controls) < non-obvious implementation choice (variance) < restraint against the visible objective (hard). AUTHZ, where all models fail, represents the extreme—restraint combined with multi-document conflict resolution.

## 9 Technical Details

- **Format:** METR Task Standard (task.yaml + repo\_snapshot + grader)
- **Runtime:** `node --test` with `node:assert/strict`, zero npm dependencies
- **Coverage:** 11 OWASP ASVS chapters + 3 OAuth 2.1 protocols
- **Hidden checks:** 51 total, 3 per environment, all deterministic and code-based
- **Calibration:** 51/51 three-state validations (17 environments × 3 states: baseline/bad/good)
- **Ablation:** 16/17 tasks probed (PRIVACY pending retest); 15/16 lose at least one hidden check without documents
- **Models tested:** Claude Sonnet 4, Claude Opus 4 (2–4 runs per task), Claude Haiku 4.5 (PRIVACY only). Model IDs: `claude-sonnet-4-20250514`, `claude-opus-4-20250514`, `claude-haiku-4-5-20251001`
- **API parameters:** Temperature: API default (1.0); max output tokens: 8,000; system prompt: “You are a careful Node.js security engineer. Return JSON only.”; response format: structured JSON with edits array
- **Bad candidates:** Hand-written with annotated failure modes, not model-generated
- **Human baseline:** 1 participant, 3 tasks, 55 minutes, 8/9 hidden checks
- **Task selection:** The 17 environments are exhaustive—every environment we built is included. No tasks were excluded based on model performance. Selection of OWASP ASVS chapters and OAuth protocols was driven by requirement density and testability, not by pilot results.

## 9.1 Limitations

**Sample size.** Most tasks have 3–4 runs per model. Pass rates are point estimates; Wilson 95% confidence intervals are wide (e.g., 0/4: [0.00, 0.71]; 1/3: [0.04, 0.79]; 2/4: [0.16, 0.84]; 3/3: [0.29, 1.00]). The Sonnet-hard and Variance tiers overlap in CI and should be interpreted as rough categories, not statistically resolved boundaries. Increasing to pass@10 would tighten these estimates substantially.

**Human baseline.**  $n = 1, 3$  tasks. The participant’s performance (8/9 hidden checks, 55 minutes) is a preliminary signal that the tasks are human-solvable, not a population estimate.

**Prompt sensitivity.** All runs use a single fixed prompt per task. We have not tested prompt wording variations. Results may differ under rephrasing, chain-of-thought elicitation, or multi-turn agent scaffolding.

**Model coverage.** Results are from three Claude models. Other model families (GPT-4, Gemini) have not been tested. The methodology is model-agnostic; the specific difficulty gradient may not transfer.

## 10 Related Work

**Specification gaming and reward hacking.** Bondarenko et al. [2] demonstrate that reasoning models (o3, DeepSeek R1) spontaneously hack benchmark environments—overwriting board state files, replacing engine binaries—in 88% of runs, while non-reasoning models game only when prompted. Wang and Huang [1] prove this is structural: under finite evaluation, effort reallocation away from non-evaluated dimensions is an equilibrium, not a failure mode. Their Proposition 2 shows that as tool count grows, evaluation coverage ratio  $K/N \rightarrow 0$  and hacking severity increases without bound. Our split-verifier architecture is a direct response: the hidden tests measure the non-contractible dimensions that visible tests leave uncovered.

**Benchmark contamination.** Prathifkumar et al. [3] show that Claude models perform 3–6× better on SWE-Bench-Verified than on structurally identical fresh benchmarks, with the gap widening when contextual cues are removed—consistent with training-data recall rather than reasoning. Our authority-ablation mode (Section 7) addresses this from the opposite direction: rather than testing whether the model has memorized the benchmark, we test whether it has memorized the *answer*. Stripping authority documents and observing hidden-check degradation confirms that correct solutions depend on document comprehension, not prior exposure.

**Operational compliance evaluation.** AgentOrca [4] evaluates constraint adherence using a dual system: natural-language prompts for agents and executable code as ground truth, across 663 tasks in five domains. Our work differs in three respects. First, our tasks feature *conflicting* authority documents where two valid rules apply and the correct resolution depends on understanding what each rule protects—not merely whether the agent follows the constraint. Second, we measure purpose recovery (understanding *why*) rather than procedural adherence (following *what*). Third, our bad-candidate calibration proves that each hidden test discriminates, which AgentOrca’s framework does not require.

**Non-monotonic scaling.** Gema et al. [5] construct tasks where extended reasoning *worsens* performance: Claude Opus 4 accuracy drops from ~100% to 85–90% on counting tasks with distractors, and Claude Sonnet 4 exhibits increased self-preservation expressions under longer reasoning traces. Our PRIVACY result (Haiku passes, Sonnet fails, Opus mostly passes) and FORGOT result (Opus

trails Sonnet) are consistent with this pattern: more capable models may overthink conflicting requirements and rationalize the wrong resolution.

## 11 Discussion

The restraint-based difficulty finding—that tasks become hard when the hidden contract opposes the visible objective—is specific enough to generate new environments systematically. Any domain where authority documents impose constraints that limit the obvious fix is a candidate: data retention rules that prevent complete deletion, backward compatibility requirements that prevent clean migration, safety invariants that prevent performance optimization, regulatory holds that prevent standard processing.

OWASP ASVS 5.0 contains 350+ verifiable requirements. OAuth 2.1 adds 50+. IEC 62443 (industrial cybersecurity) adds a third source. Each requirement is convertible to an environment using the methodology described here. The standards are versioned, so new requirements produce new tasks as specifications evolve.

The methodology is not security-specific. The same hidden-contract pattern applies wherever expert-authored documents govern system behavior: compliance frameworks, medical protocols, engineering specifications, financial regulations. Security is the proof of concept.

## Availability

The full suite—17 environments, all calibration candidates, hidden tests, authority documents, and run harness—is publicly available at <https://github.com/cv700/asvs-security-eval> under CC BY-NC 4.0.

Cooper Veit  
Ashiba Research  
coopveit@gmail.com

## References

- [1] J. Wang and J. Huang. Reward hacking as equilibrium under finite evaluation. *arXiv preprint arXiv:2603.28063*, March 2026.
- [2] A. Bondarenko, D. Volk, D. Volkov, and J. Ladish. Demonstrating specification gaming in reasoning models. *arXiv preprint arXiv:2502.13295*, February 2025.
- [3] T. Prathifkumar, N. S. Mathews, and M. Nagappan. Does SWE-Bench-Verified test agent ability or model memory? *arXiv preprint arXiv:2512.10218*, December 2025.
- [4] Z. Li, S. Huang, J. Wang, N. Zhang, et al. AgentOrca: A dual-system framework to evaluate language agents on operational routine and constraint adherence. *arXiv preprint arXiv:2503.08669*, March 2025.
- [5] A. P. Gema, J. Hagele, Y. Chen, et al. Inverse scaling in test-time compute. *Transactions on Machine Learning Research*, December 2025.

- [6] T. Kwa, B. West, J. Becker, et al. Measuring AI ability to complete long software tasks. *arXiv preprint arXiv:2503.14499*, March 2025.
- [7] T. Korbak, M. Carroll, B. Baker, and I. Kivlichan. Reasoning models struggle to control their chains of thought, and that's good. *arXiv preprint arXiv:2603.05706*, March 2026.



## Ashiba Research

Ashiba Research is a fresh-formed data provider building datasets for the messy real world. Cooper Veit combines consultant-GTM operational chops (two years in consulting, \$150M founding GTM exit) and editorial discipline in problem selection and qualification. Ashiba is focused on the development of environments in which highly compressed and constrained systems can compose into reliable systems.

Ashiba can be a dataset vendor or scout. We have strong GTM discipline and sightlines to standards bodies, diagnosticians and workers in neglected corners of the economy, struggling small businesses, and academic labs and tech transfer offices.

However, we believe that the most direct path to useful data runs through evals and systems that we design ourselves. Ashiba is currently developing three research programs with the aim of converging on one after making contact with the frontier:

1. **Prime Standard** — Something like Antithesis for messy technical standards, protocols, and codebases. We editorially select standardized markets such as cyber, determine where the bugs are in existing protocols, and imbue agents with useful compressed ontologies or crystalline problem knowledge that drives superior performance when dealing with those standards and others. These methods should enable the production of benchmarks, evals, environments, and standards across high-TAM existing and developing markets starting with cyber. We sell evaluation and training environments to frontier AI labs; experts earn royalties on the reasoning embedded in each. *The preceding report is a sample of this work.*

2. **Ashiba Compute** — Kernel contracts and silent-correctness conformance for the fragmented AI silicon stack. Every ML kernel ships with an implicit contract about what it computes—numerical tolerance, determinism, shape limits, composition semantics—but the contract is almost never written down. When two kernels disagree (a matmul on AMD and NVIDIA producing different gradients, Huawei Ascend silently downcasting an accumulator, an out-of-bounds access returning zero on one stack and garbage on another), there is no formal artifact to arbitrate the dispute. We codify a 12-class taxonomy of silent correctness failures into a formal kernel contract language, publish it as an open reference under Ashiba, and operate a continuously-updated conformance dataset measuring which hardware/software combinations meet which contracts across AMD ROCm, Intel Gaudi, Huawei Ascend, AWS Trainium, and Triton backends. Closer to Chainalysis than SemiAnalysis in shape—forensic data about vendors who will not self-report. Sibling infrastructure to Prime Standard and ProbSpec: same licensing pattern, third domain. We sell a subscription conformance dataset to frontier labs and hyperscalers diversifying off CUDA; certification royalty flow compounds once the contract suite is cited by a standards body.
  
3. **ProbSpec** — A licensable ontology and certification scheme for the interop layer that governs legacy industrial systems. We instrument retiring tradespeople through community-college partnerships to capture labeled fault→diagnosis→fix episodes, codify them into 12 typed failure classes—alongside supplier reliability schemas, recovery playbooks, and compatibility matrices—and seed the resulting technical report inside the IEC 62443 industrial cybersecurity working group. The endgame is a brownfield-interoperability sub-scheme under ISASecure, licensed through accredited certification bodies (exida, TÜV, DNV) to integrators, OEMs, and asset owners, anchored by EU Machinery Regulation obligations taking effect January 2027. Sibling infrastructure to Prime Standard—same licensing pattern, applied to the industrial domain. We sell expert-trace datasets to frontier labs, supplier-reliability intelligence to insurance carriers and OEM procurement, and conformance royalties once the ISASecure sub-scheme launches.

Parallel closers, three distinct revenue shapes:

Program	Primary buyer	Product	Moat
Prime Standard	Frontier AI labs	Eval & training envs	Methodology + expert royalties
Ashiba Compute	Labs + hyperscalers	Conformance data (sub.)	Kernel contract reference
ProbSpec	Labs + insurers + OEMs	3 data products + royalty	Failure ontology + ISASecure

Sightlines to cofounders and early teams.