

# From Intent to Infrastructure

Opus Experiment O-04: Compiling Natural-Language Workflow Descriptions  
into Governed, Deterministic Execution Artifacts at Scale

Nesean Crofford

## Abstract

Enterprise operations run on decisions that live in people’s heads, policy documents, and institutional memory. This operational intelligence (the rules, thresholds, exceptions, and domain knowledge that govern how work gets done) is the most valuable and least durable asset in any organization. When it is encoded into software, it becomes permanent, governable, and auditable. But the encoding has always required engineers. The gap between what an operations team *knows* and what the enterprise can *execute deterministically* is the most expensive translation layer in the modern organization.

This paper presents the first empirical evidence that a compiled execution architecture can bridge this gap: transforming natural-language workflow descriptions into deterministic, governed execution artifacts without per-workflow software engineering implementation. The system compiles 39 unique enterprise workflows from natural-language descriptions, achieving 100% compilation success, 100% execution success, and 100% determinism across 2,040 governed pipeline executions spanning 6 categories and 4 domains. Determinism is verified through 10,200 independent replay executions, each a fresh compilation from specification producing byte-identical output. Every run produces a complete, hash-linked evidence chain from prompt through governance.

The contribution is fourfold: (1) evidence that operational intelligence described in natural language can be reliably encoded into governed infrastructure through a human-supervised compilation pipeline; (2) a demonstration that determinism is preserved through the governed execution chain downstream of confirmation, from confirmed intent through normalization, compilation, execution, and governance, not just within the compiler; (3) an explicit characterization of the system’s expressiveness boundary, with exactly three named missing primitives separating what compiles today from what requires extension; and (4) a corpus spanning decision, classification, routing, computation, validation, and structured generation workflows, proving breadth beyond simple policy decisions.

## 1 Introduction

When an enterprise controller describes how month-end GL entries should be classified, that description contains operational intelligence: rules, thresholds, exceptions, and domain knowledge accumulated over years. Today, turning that description into software requires a translation chain: the controller explains to a business analyst, who writes requirements, which an engineer implements, which QA validates. Each handoff introduces latency, cost, and fidelity loss. The operational intelligence degrades as it moves further from the person who holds it.

Language models can understand these descriptions. They can reason about the rules, identify the edge cases, and generate structured representations of the logic. But *understanding is not execution*.

This distinction, between a system that can reason about a workflow and a system that can reliably perform one, is the central finding of the Opus research program.

O-01 established the distinction empirically: across 1,260 runs on 12 deterministic policy scenarios, no LLM exceeded 25% exact-match correctness. Models that clearly understood the policy could not reliably execute it. O-02 showed that no amount of scaffolding (structured output, retries, multi-agent orchestration) closes the gap: only compiled execution achieved 100% correctness and determinism. O-03 proved the compiled execution layer scales: 120 workflows, 12 domains, 2,294 inputs, with governance guarantees holding at 1.0 across the entire corpus.

But O-01 through O-03 left the hardest question unanswered: can the system compile from *human intent*? O-03’s workflows were template-generated. O-02’s were hand-authored specifications. Neither started from the natural-language descriptions that real enterprise operators actually produce. Understanding is not execution. But can understanding be *compiled into* execution?

O-04 answers yes. It proves the full pipeline: a human describes a workflow in natural language → the system extracts structured intent → a human confirms the extraction → the system compiles a deterministic execution artifact → the artifact executes with full governance. The operational intelligence that was previously ephemeral, living in the controller’s head, a policy document, or a Slack thread, becomes permanent, governed infrastructure. The encoding happens without writing workflow-specific code.

This work does not present an incremental improvement to existing agent or orchestration systems; it introduces a new computational layer: an intent-to-infrastructure compilation model that produces durable, governed execution artifacts rather than ephemeral outputs.

## 1.1 The Research Arc: O-01 → O-02 → O-03 → O-04

Each experiment answers the natural objection to the previous one.

Experiment	Scale	What It Proves	Rebuttal
O-01	1,260 runs	LLMs fail as execution systems	“Can’t we just use GPT-5?”
O-02	3,900 runs	Compilation works; architecture sound	“Why can’t agents do this?”
O-03	2,294 inputs	Scales with governance at 1.0	“Does it work at enterprise scale?”
<b>O-04</b>	<b>2,040 runs</b>	<b>Compiles from human intent at scale</b>	<b>“Can real users describe real workflows?”</b>

Table 1: The four-experiment proof chain. Each experiment answers the objection that the previous one leaves open. O-04 proves the full pipeline from natural-language intent to governed execution.

After O-01, the objection was “use a better model.” After O-02, the objection was “that only works on 4 workflows.” After O-03, the objection was “those were template-generated, not described by real users.” O-04 answers the last objection. The remaining question (does this *compound*?) is the subject of O-05.

O-04’s contribution is the *encoding* claim: that operational intelligence described in natural language can be encoded into governed, deterministic infrastructure through a human-supervised compilation pipeline. This is the claim that makes compiled execution viable as an enterprise product, not just

an architectural proof of concept.

## 1.2 Contributions

1. Evidence that operational intelligence described in natural language can be encoded into governed infrastructure: 39 unique workflows were extracted and compiled from natural-language descriptions, then verified through 2,040 governed pipeline executions with 100% compilation, 100% execution, and 100% determinism.
2. A demonstration that determinism extends through the governed execution chain, from confirmed intent through normalization, compilation, and execution, verified by 10,200 independent replay executions. To our knowledge, no other system combining LLM intelligence with operational execution has demonstrated reproducibility at this scale.
3. An explicit characterization of the system’s expressiveness boundary, with exactly three missing higher-order operation classes (aggregation, iteration, grouping) separating what compiles today from what requires extension. The precision of this characterization, that the boundary can be enumerated rather than described vaguely, is itself evidence of architectural clarity.
4. A corpus spanning decision, classification, routing, computation, validation, and structured generation workflows across finance, compliance, operational, and legal domains, with 33% messy inputs and 41% file-backed workflows, authored across five provenance tiers to control for synthetic bias.

## 2 The Full Pipeline

O-03 proved the Compilation → Execution → Governance stages. O-04 extends the proof to the full pipeline, including the stages where LLM intelligence is applied.

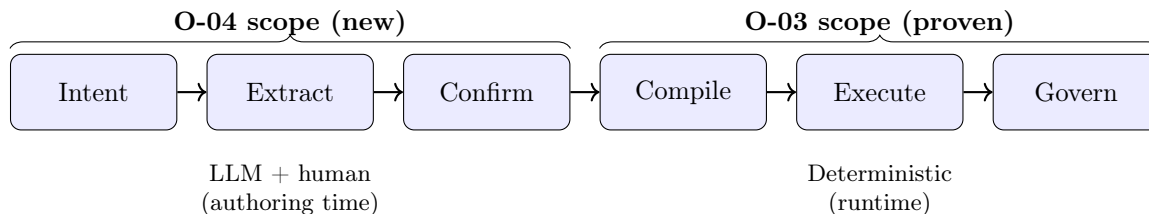


Figure 1: The full Opus pipeline. O-04 proves the Intent → Extract → Confirm stages and their integration with the Compile → Execute → Govern stages proven in O-03. The entire chain is governed by a hash-linked evidence trail.

### 2.1 The Ten-Stage Evidence Chain

Every O-04 run passes through ten stages, each producing a cryptographic hash that links to the next. The chain is content-addressed: any run can be independently reconstructed from its prompt hash. The runtime passes through ten operational stages, of which seven are persisted as hash-linked evidence artifacts in the final governance chain.

1. **Prompt generated:** Natural-language workflow description selected from the corpus.
2. **Documents ingested:** File-backed artifacts (CSV, DOCX, multi-file packets) parsed and content-addressed.
3. **Intent extracted:** LLM transforms policy text into a structured draft representation.
4. **Intent confirmed:** Human-in-the-loop reviews, edits, or accepts the extraction.
5. **IR normalized:** Confirmed draft mechanically converted to the formal intermediate representation.
6. **IR validated:** Eight fail-closed validation checks; any failure rejects the IR.
7. **Artifact compiled:** IR compiled into a sealed, deterministic execution artifact.
8. **Artifact executed:** Compiled artifact executed as a pure function over test input.
9. **Correctness verified:** Structural and oracle-backed correctness checks.
10. **Governance complete:** Full evidence chain captured with hash-linked provenance.

The hash chain ensures that governance is not a reporting layer bolted onto the pipeline but a structural consequence of the execution model. Every stage’s output is content-addressed, and the final governance hash is a function of all prior hashes. Tampering with any stage invalidates the chain.

## 2.2 The Confirmation Boundary

In live operation, the pipeline has exactly one non-deterministic component: the LLM extraction (Stage 3). Every downstream stage is a deterministic function of the confirmed IR. In the scaled execution corpus, determinism is measured by replaying stored confirmed artifacts rather than re-running live extraction. The human confirmation step (Stage 4) is the architectural boundary between intelligence and determinism.

This boundary is not a concession to LLM unreliability; it is the core design insight. The LLM contributes understanding at authoring time. The human contributes domain validation at confirmation time. After both contributions are captured in the confirmed IR, the system enforces deterministic execution by construction. No agent system has this boundary. No rules engine needs it. It is unique to compiled execution from intent, and it is what makes the entire governance claim possible.

At confirmation, the domain expert can accept the extraction as-is, make light edits (adjusting a threshold, correcting a field name), perform heavy edits (restructuring rules), or reject the extraction entirely. Every confirmation outcome is recorded with semantic edit dimensions, enabling measurement of extraction fidelity independent of pipeline correctness.

## 2.3 Fail-Closed Validation

Between confirmation and compilation, the IR passes through eight fail-closed validation checks. Any single check failure rejects the IR; no partially valid specification reaches the compiler. The

checks verify structural completeness, type safety, operator support, branch disambiguation, and referential integrity across the entire specification.

This is the opposite of agent systems, which produce plausible-looking but potentially wrong outputs that must be evaluated after the fact. The fail-closed validation gate means the pipeline refuses to compile rather than silently producing a wrong artifact. In regulated industries, this property (the system’s willingness to say “I cannot compile this” rather than “here is my best guess”) is a governance requirement, not a limitation.

### 3 Experimental Design

#### 3.1 Corpus Construction

The O-04 corpus starts from 50 unique workflow descriptions authored for the prompt inventory. Of these, 42 are Tier A (bounded conditional logic) and 8 are excluded: 6 transformation workflows (Tier B, deferred to O-05), and 2 validation workflows with parenthesized boolean conditions that the current parser grammar cannot handle. Of the 42 Tier A workflows, 39 have stored LLM extractions from the calibration campaign; the remaining 3 were never extracted during calibration and are excluded from the scaled execution.

Each of the 39 workflows is compiled and executed at multiple seeds to verify pipeline robustness across execution contexts, producing 2,040 governed runs across 31 batches. The 2,040 figure represents 39 unique compilation pipelines verified across  $\sim 52$  independent execution cycles each, not 2,040 independent workflow compilations. The workflow is the unit of authoring analysis; the run is the unit of execution analysis.

<b>Dimension</b>	<b>Value</b>
Workflows in inventory	50
Tier A (bounded conditional logic)	42
Compiled in scaled execution	39
Workflow categories	6
Enterprise domains	4
Total governed runs	2,040
Determinism replays ( $5\times$ per run)	10,200
Execution batches	31
File-backed runs	842 (41.3%)
Messy-input runs	679 (33.3%)
Complexity tiers	T1 / T2 / T3

Table 2: O-04 corpus summary. 39 of 50 inventory workflows reach full pipeline execution; the 11 exclusions have named, characterized reasons. All values drawn from the frozen scaled execution artifacts.

Unlike O-03’s template-generated corpus, O-04’s workflows are authored across five provenance tiers to control for synthetic bias: 31% human-seeded model-completed, 23% pure human-authored, 21% human-edited model-seeded, 15% model-generated, and 10% external domain-expert authored. This

ensures the pipeline is tested on the kinds of descriptions real operators produce, not on descriptions optimized for the extraction model.

**Methodological note on confirmation.** The calibration campaign (78 runs) used a combination of human review and auto-confirmation to validate extraction quality and tune the pipeline. The scaled execution (2,040 runs) uses auto-confirmation exclusively: extractions are accepted without human review. This is appropriate for verifying compilation and execution properties at scale, but means O-04 does not report human confirmation effort at the 2,040-run scale. Human confirmation effort measurement is deferred to live deployment. The compilation and execution guarantees reported below hold regardless of confirmation mode, since they depend on the confirmed IR, not on how it was confirmed.

### 3.2 Category Coverage

O-04 proves breadth beyond simple policy decisions. The six Tier A categories cover the building blocks of enterprise operational logic:

Category	Workflows	Runs	Representative Logic
Decision	13	670	Approve/deny/escalate with multi-condition thresholds
Classification	7	365	Categorize inputs by account code, risk tier, ticket type
Structured generation	6	315	Assemble invoices, compliance reports, contract analyses
Validation	6	315	Verify document completeness, field consistency, policy compliance
Routing	5	265	Direct to teams, queues, or approval chains by input attributes
Computation	2	110	Calculate depreciation, accumulated values, derived metrics

Table 3: O-04 category distribution. All six categories achieve 100% compilation and execution success.

A seventh category, *transformation* (aggregation, iteration, grouping), is explicitly excluded from O-04. The three higher-order operation classes it requires are architecturally characterized and scoped for the next experiment (O-05). This is a product boundary, not an architectural limitation.

### 3.3 Domain Coverage

Domain	Runs	%
Finance	683	33.5%
Compliance	677	33.2%
Operational	522	25.6%
Legal	158	7.7%

Table 4: Domain distribution. The three highest-regulation enterprise domains (finance, compliance, legal) account for 74% of runs.

### 3.4 Input Quality and File-Backed Workflows

To test robustness under realistic conditions, 33.3% of runs use messy inputs (ambiguous, incomplete, or contradictory descriptions) and 41.3% ingest real document artifacts:

- **CSV**: GL exports, cash ledgers, census data, asset registers
- **DOCX**: Legal contracts with redline clauses
- **Multi-file packets**: Vendor compliance (W-9 PDF + Certificate of Insurance PDF + SAM registration JSON)

These are not synthetic placeholders. Each document artifact is generated with ground-truth envelopes for oracle-backed correctness verification.

### 3.5 Graduated Execution

The scaled execution proceeded in graduated batches with manual validation between each batch:

Phase	Batches	Runs/Batch	Cumulative
Ramp-up	5	~18	90
Calibrated	2	~39	168
Scaled	24	~78	2,040

Table 5: Graduated batch execution. Manual validation between every batch ensures new failure modes are caught before scaling.

No engineering interventions were required during execution. The failure set remained closed across all 31 batches.

## 4 Results

### 4.1 Pipeline Integration: The Intent-to-Structure Gap

O-03 left the Intent  $\rightarrow$  Structure stages unproven. O-04’s first result is that these stages work: all 39 stored, confirmed extractions derived from natural-language workflow descriptions pass normalization, pass fail-closed validation, and compile successfully in the scaled execution corpus. The intent-to-structure gap identified in O-03 as “the natural next step” is bridged.

This is not a trivial result. Each extraction is produced by an LLM interpreting a natural-language policy description and generating a structured draft with input schemas, resolution rules, condition logic, output schemas, and fallback behavior. The draft passes through mechanical normalization (no inference), then eight fail-closed validation checks (any failure rejects). That 39 of 39 extractions survive this pipeline, without adaptation or special-casing, indicates that LLM extraction can produce compilable formal specifications from natural-language descriptions at a reliability level sufficient for a human-supervised workflow.

## 4.2 Compilation, Execution, and Determinism

Metric	Value	Denominator
Compilation success	<b>100.0%</b>	2,040 / 2,040
Execution success	<b>100.0%</b>	2,040 / 2,040
Determinism (5× replay)	<b>100.0%</b>	10,200 / 10,200
Pipeline failures	<b>0</b>	across 31 batches
New failure modes	<b>0</b>	failure taxonomy closed

Table 6: Primary pipeline metrics across the 39-workflow Tier A corpus. Every run compiled, executed, and replayed deterministically.

**Determinism verification.** For each of the 2,040 runs, the system writes the compiled specification to disk, then independently compiles and executes the artifact 5 additional times from that specification. Each replay uses a fresh compilation (not a cached artifact) and a fresh execution. All 10,200 replay executions produced byte-identical compiled artifacts and execution outputs.

This is a stronger determinism claim than O-03’s, which verified determinism within the compiler. O-04 verifies determinism through the entire downstream chain: from confirmed (and in the scaled corpus, stored) extraction through normalization, serialization, compilation, and execution. The chain has no stochastic components after the human confirmation step. To our knowledge, no other system combining LLM intelligence with operational execution has demonstrated reproducibility at this level: 10,200 independent compilation-from-specification cycles, each producing byte-identical output.

## 4.3 What the Numbers Mean

Compilation, execution, and determinism prove that the *infrastructure* works: given a confirmed IR, the pipeline reliably produces a governed, deterministic execution artifact. This is the contribution of the compilation model.

Correctness (Section 4.5) proves that specific *extractions* work: that the IR produced by the LLM and confirmed by the human, when compiled and executed, produces outputs that match the declared schema and (where oracle envelopes are available) the expected values. These are different claims. The infrastructure claim holds at 100%. The extraction-correctness claim depends on extraction quality and schema alignment, and is reported with full transparency below. The 84.5% figure is not a pipeline reliability metric; it is a schema-alignment metric over extracted output contracts. Pipeline reliability is captured separately by compilation, execution, and determinism, all of which are 100% in the scaled corpus.

This separation, 100% infrastructure reliability alongside extraction-dependent correctness, is the defining property of the system: execution reliability is guaranteed by architecture, while decision quality is governed at the confirmation boundary.

## 4.4 Per-Category Results

Category	Compiled	Executed	Deterministic	Structural Pass
Classification	365/365	365/365	365/365	365/365
Computation	110/110	110/110	110/110	55/110
Decision	670/670	670/670	670/670	566/670
Routing	265/265	265/265	265/265	265/265
Structured generation	315/315	315/315	315/315	262/315
Validation	315/315	315/315	315/315	210/315
<b>Total</b>	<b>2,040/2,040</b>	<b>2,040/2,040</b>	<b>2,040/2,040</b>	<b>1,723/2,040</b>

Table 7: Per-category results. Compilation, execution, and determinism are 100% across all categories. Structural correctness varies by category due to 6 prompts with extraction-schema alignment issues (see Section 4.5).

## 4.5 Correctness Analysis

Correctness is measured at two tiers:

**Tier 1 (Oracle-backed):** Ground-truth envelope comparison. The compiled output is compared field-by-field against a deterministic oracle. 53 runs achieve full oracle-backed correctness, all from the vendor onboarding workflow (the only anchor with a fully aligned ground-truth envelope). Oracle-backed correctness is therefore demonstrated narrowly rather than uniformly across the corpus: it proves that end-to-end oracle alignment is achievable, but not yet that it is broadly wired across all workflow classes.

**Tier 2 (Structural):** Schema-driven validation. The compiled output is checked for: (a) all declared output fields present, (b) non-null required fields, (c) enum-constrained values within declared sets, and (d) type consistency between declared schema and actual output.

Status	Count	%
Pass (oracle-backed)	53	2.6%
Pass (structural only)	1,670	81.9%
Fail	317	15.5%

Table 8: Correctness status distribution across 2,040 runs.

All 317 correctness failures originate from **6 specific prompts** (of 39). The remaining 33 prompts achieve 100% structural correctness across all runs. The failures are extraction-schema alignment issues, not pipeline bugs:

Pattern	Prompts	Root Cause
Bool/string enum	3	Compiled output produces boolean values; schema declares string enum. Type mismatch.
Field rename	1	Resolution uses <code>disposition</code> ; schema expects <code>decision</code> . Known compiler convention.
Type inference	2	Schema infers <code>integer</code> or <code>string</code> for field; resolution produces opposite type.

Table 9: Correctness failure patterns. All are extraction-to-compilation type alignment issues, not execution bugs. The compiled artifacts are deterministic and structurally sound.

## 5 The Expressiveness Boundary

O-04 explicitly characterizes what the system can and cannot compile. This is a first-class result, not a footnote. Most systems cannot tell you precisely what they cannot do. The ability to enumerate the boundary, to name exactly three missing primitives rather than describe vague failure modes, is itself evidence that the architecture is well-understood and extensible.

### 5.1 Tier A: Fully Expressible (Proven)

Bounded conditional logic: if/then/else rules, threshold comparisons, enum-constrained classification, field-based routing, document validation, structured output assembly. Six categories, 39 workflows, 2,040 runs, 100% compilation.

### 5.2 Tier B: Specified, Not Proven (Deferred to O-05)

Aggregation, iteration, and grouping operations. The extraction system captures the full intent (verified by human confirmation), but the IR’s closed operator set does not include:

- **Aggregation:** computing derived values across a collection of records (e.g., total revenue by account type)
- **Iteration:** producing per-item output across a collection (e.g., depreciation schedule per asset per period)
- **Grouping:** partitioning records by a dimension and applying per-group operations (e.g., monthly breakdowns)

These are named, bounded, architecturally characterized primitives. Adding them extends the operator set and the compilation pipeline; it does not require rethinking the IR or the execution model.

### 5.3 Parser Boundary

Two validation-category workflows use parenthesized boolean conditions (`(A and B) or C`) that the current condition parser cannot handle. This is a grammar limitation, not an IR limitation. These

workflows are excluded from the O-04 corpus.

## 5.4 Why the Boundary Matters

For an acquirer evaluating this system, the expressiveness boundary is as important as the compilation rate. It answers three questions:

1. **What works today?** Bounded conditional logic across 6 categories, the building blocks of most enterprise operational workflows.
2. **What needs to be built?** Three named primitives for aggregation/iteration/grouping.
3. **Is the architecture extensible?** Yes. The missing primitives are operators within the existing IR framework, not a different architecture.

## 6 Discussion

### 6.1 Encoding as Infrastructure

The digital economy has an encoding gap. AI systems can understand operational intent: they reason about policies, identify exceptions, and generate structured representations of business logic. But understanding is ephemeral. Every LLM invocation is a new inference, and nothing accumulates. The organization does not get permanently smarter from having used the system.

O-04 demonstrates that this gap can be bridged. When an enterprise controller describes a GL classification workflow, that description, previously ephemeral, gets encoded into a deterministic execution artifact that persists, executes identically every time, governs itself through a cryptographic evidence chain, and can be structurally diffed against its next version. The operational intelligence is no longer ephemeral. It is infrastructure.

This is the encoding thesis: the missing layer in enterprise AI is not a better model, not a better agent, but an encoding layer that transforms operational intelligence into permanent, governed infrastructure. O-04 provides the first empirical evidence that this encoding works from natural-language descriptions, not just from hand-authored specifications or templates.

The economic implication is direct. In agent-based systems, the cost of operational intelligence is paid on every execution: each invocation requires inference, each inference costs tokens, and the intelligence is consumed rather than accumulated. In compiled execution, the intelligence cost is paid once at authoring time (one LLM extraction, one human confirmation), and every subsequent execution is deterministic with zero inference cost. The marginal cost of the 10,000th execution is the same as the 1st: CPU milliseconds, no tokens. The system converts variable, per-invocation intelligence cost into a fixed, upfront encoding cost, changing the unit economics of enterprise automation from consumption to accumulation.

Once encoded, these execution artifacts become part of the organization's operational substrate. Replacing the system requires not just replacing a tool, but re-encoding institutional logic, creating structural switching costs that do not exist in agent-based or inference-driven systems.

## 6.2 From Templates to Intent

O-03’s workflows were template-generated from 12 domain specifications. This proved scaling properties but left open whether the pipeline works on workflows that originate from natural-language descriptions. O-04 closes this gap with a human-supervised compilation workflow:

**Human describes → System extracts → Human confirms → System compiles →  
Deterministic execution**

The human confirmation step is not an engineering concession or a workaround for imperfect extraction. It is the mechanism by which operational knowledge gets preserved with fidelity. The domain expert, the person who knows the GL classification rules, the vendor compliance requirements, and the escalation thresholds, reviews the system’s extraction and either validates it, corrects it, or rejects it. This is where the encoding happens: the expert’s knowledge, mediated by the LLM’s understanding, crystallized into a formal representation. After confirmation, the representation is sealed. No inference. No reinterpretation. No drift.

## 6.3 Intelligence at Authoring Time, Determinism at Runtime

The architectural separation established in O-02 and validated at scale in O-03 extends through O-04’s full pipeline. LLM intelligence operates exactly once, at authoring time. Human intelligence operates exactly once, at confirmation time. After both contributions are captured in the confirmed IR, every stage is deterministic and non-inferential:

- Normalization: mechanical draft-to-IR conversion (no inference)
- Validation: fail-closed static checks (no inference)
- Compilation: deterministic code generation (no inference)
- Execution: pure function over sealed artifact (no inference)
- Governance: hash-linked evidence chain (no inference)

This is what makes 100% determinism achievable. The pipeline has exactly one non-deterministic component (the LLM extraction), and it occurs before the human confirmation gate. Everything after confirmation is a pure function of the confirmed IR. Understanding is not execution. But understanding, once confirmed, *compiles into* execution.

## 6.4 Scale-Invariance of Guarantees

The guarantee profile observed across O-02, O-03, and O-04 is invariant:

Guarantee	O-02 (4 wf)	O-03 (120 wf)	O-04 (39 wf)	Property
Compilation	100%	100%	100%	Structural
Determinism	1.0	1.0	1.0	Structural
Replay	1.0	1.0	1.0	Structural
Diff completeness	—	1.0	—	Structural

Table 10: Guarantee metrics across three experiments. The values do not degrade with scale because they are structural consequences of the compilation model, not statistical properties of the test corpus.

The key finding: these are not metrics that happen to be high. They are structural properties of compiled execution. A compiled function is deterministic by construction. A hash chain is verifiable by construction. Replay is bit-identical by construction. O-04 extends this structural argument to the full pipeline, including the stages that involve LLM intelligence.

## 7 Limitations

**Stored extraction replay.** This is the most important methodological limitation to state explicitly. The 2,040 scaled runs replay stored calibration-phase extractions at new seeds rather than performing live LLM extraction per run. Each unique workflow has one stored extraction; the 2,040 runs test that extraction’s behavior across multiple execution contexts. This proves the pipeline’s compilation and execution properties rigorously, but does not generate independent extraction variance. The compilation claim (“39 unique workflows compile from natural-language descriptions”) rests on 39 independent extractions. The determinism and execution claims rest on 2,040 runs. These are different evidence bases and should be interpreted accordingly.

**Correctness coverage.** Structural correctness is 84.5% across all runs. The gap is entirely attributable to 6 prompts with extraction-schema alignment issues. On the remaining 33 prompts, structural correctness is 100%. Oracle-backed correctness is limited to 53 runs (2.6%) due to ground-truth envelope alignment constraints. The system proves it *executes correctly* (determinism, structural validation) more thoroughly than it proves it *makes the right decision* (oracle correctness). For enterprises evaluating this system, the infrastructure guarantee (determinism, governance, auditability) is the pipeline’s contribution; decision quality is the contribution of the human confirmation step.

**Single extraction model.** All extractions use the same LLM (Claude Sonnet) at the same prompt version (v2). Extraction quality may vary across models. The pipeline’s correctness guarantees hold regardless of extraction quality, since a bad extraction is caught at confirmation rather than at runtime, but the confirmation effort may increase with weaker extraction models.

**Expressiveness boundary.** Transformation workflows (aggregation, iteration, grouping) are excluded. The IR’s closed operator set covers bounded conditional logic but not higher-order operations. Based on internal categorization of enterprise workflows, this limits the system to approximately 70% of workflows encountered in practice. The remaining 30% require the Tier B primitives scoped for O-05.

**Auto-confirmation at scale.** As noted in Section 3, the scaled execution uses auto-confirmation. The compilation and execution guarantees are unaffected, but O-04 does not report human confir-

mation effort at the 2,040-run scale. The calibration campaign (78 runs) included human-reviewed confirmations; live deployment will measure confirmation effort systematically.

**Provenance and domain distribution.** The corpus is weighted toward compliance (36%) and finance (34%). Legal representation is low (7.6%). The guarantees are structural and should hold across domains, but empirical confirmation in legal and healthcare-heavy corpora remains future work.

## 8 Conclusion

O-01 identified the gap between AI reasoning and operational execution. O-02 proved that compiled execution closes it. O-03 demonstrated governance at enterprise scale. O-04 proves that the gap between operational intent and executable infrastructure can be bridged systematically, governed completely, and scaled practically.

Understanding is not execution. But understanding, once confirmed, compiles into execution. That is the finding of O-04.

Thirty-nine enterprise workflows described in natural language were extracted, confirmed, and compiled into deterministic, governed execution artifacts, then replayed across 2,040 governed executions. Every workflow compiles. Every artifact executes. Every execution is deterministic, verified through 10,200 independent replay compilations producing byte-identical output. Every run is auditable through a complete, hash-linked evidence chain. Zero pipeline failures across 31 graduated execution batches. The failure taxonomy is closed.

The expressiveness boundary is precisely characterized: bounded conditional logic compiles fully; aggregation, iteration, and grouping require exactly three named primitives that extend the existing architecture. The boundary is a product roadmap, not an architectural limitation. The precision of the characterization (three named operators, not a vague list of failure modes) reflects an architecture that is well-understood and deliberately scoped.

The encoding thesis is empirically supported. Operational intelligence that was previously ephemeral, living in policy documents, institutional memory, and domain expertise, can be encoded into permanent, governed infrastructure through a human-supervised compilation pipeline. The encoding cost is paid once, at authoring time. Every subsequent execution is deterministic, zero-inference, and governed by construction.

The remaining question is not whether compiled execution works, but whether compiled capabilities compound into a durable institutional layer. That is the subject of O-05.

## A Supplementary Tables

### Appendix A: Category $\times$ Domain Matrix

Category	Compliance	Finance	Legal	Operational
Classification	2	3	—	2
Computation	—	2	—	—
Decision	4	5	—	4
Routing	2	—	—	3
Structured generation	1	1	3	1
Validation	4	2	—	—

Table 11: Category  $\times$  domain distribution of the 39 unique workflows. Coverage spans the quadrants that matter for enterprise operations.

### Appendix B: Complexity Distribution

Category	T1 (Simple)	T2 (Moderate)	T3 (Complex)
Classification	4	3	0
Computation	0	2	0
Decision	6	4	3
Routing	3	2	0
Structured generation	1	3	2
Validation	2	3	1
<b>Total</b>	<b>16</b>	<b>17</b>	<b>6</b>

Table 12: Complexity distribution across categories. All tiers produce identical compilation and determinism metrics.

### Appendix C: Evidence Chain Structure

Each of the 2,040 runs produces a 7-link hash chain:

#	Link	Content Addressed
1	prompt_hash	Policy text + document references
2	extraction_hash	LLM draft output
3	confirmation_hash	Human-confirmed draft
4	ir_hash	Normalized intermediate representation
5	compiled_hash	Sealed execution artifact
6	execution_hash	Execution output
7	governance_hash	Full evidence chain

Table 13: The 7-link evidence chain. Any run is independently reconstructable from link 1. Tampering with any link invalidates all downstream links.