

O-Lang: Runtime Governance for Safe AI Execution in Regulated Domains

Author: Olalekan Ogundipe **Affiliation:** O-Lang Central **Date:** February 2026 **License:** Apache 2.0

Abstract

Autonomous AI agents pose structural unsafety risks in regulated domains due to unbounded execution authority, opaque decision-making, and silent data fabrication. We present **O-Lang Protocol**, an open semantic governance protocol that enforces a runtime boundary separating AI intent from execution authority. By mediating every capability invocation against explicit policy, O-Lang ensures workflows remain auditable, deterministic, and compliant—even when using non-deterministic components like LLMs. We demonstrate its application in healthcare resource allocation, show how it prevents silent failures, and contrast it with conventional orchestration frameworks. O-Lang is not a developer tool—it is infrastructure for governable AI.

1. Introduction

The rise of agentic AI has exposed a critical gap: **autonomy without accountability**. Systems like LangChain enable developers to compose multi-step workflows using LLMs and external tools. However, these systems operate *inside application code*, where:

- Developer authority = execution authority
- Tools inherit full system permissions implicitly
- Failures are masked by "graceful degradation" (e.g., hardcoded fallback values)
- No verifiable audit trail exists for regulatory compliance

This creates unacceptable risk in domains like healthcare, finance, and public services—where AI decisions directly impact human dignity, legal liability, and safety.

The O-Lang Protocol addresses this by moving governance *outside* application logic into a **runtime-enforced substrate**. It does not replace LangChain; it governs it. The core insight is simple but powerful:

AI may propose actions—but only the kernel may permit them.

2. Architectural Model

2.1 Core Components

The O-Lang Protocol consists of three normative components:

Component	Role	Security Property
Workflow	Declarative intent (what should happen)	Immutable, human-readable, versionable
Kernel	Runtime enforcer (whether it may happen)	Mediates all capability invocations
Resolver	External capability (how it is implemented)	Conformance-tested, never trusted

The kernel is the **governance boundary**. It parses workflows, validates resolvers, tracks symbol lifecycles, and emits cryptographically-verifiable audit traces.

2.2 Key Guarantees

The O-Lang Protocol provides four non-negotiable properties:

1. **Intent-Execution Separation** — Workflows declare intent; the kernel enforces execution policy. No resolver executes without explicit allow-listing.
2. **Fail-Fast Semantics** — If a symbol is undefined (e.g., database timeout), the kernel halts with `UNRESOLVED_PLACEHOLDERS` —never fabricating data.
3. **Deterministic Auditability** — Identical inputs → identical execution traces across all compliant kernels. Content may vary (e.g., LLM text), but structure does not.
4. **Explicit Failure Handling** — Retry logic belongs to workflows—not hidden in resolvers. All attempts are logged: `attempt_1: UNDEFINED`, `attempt_2: DEFINED`.

These properties ensure O-Lang Protocol systems are **certifiable**, not just functional. The principle of Deterministic Auditability draws directly from Lee's foundational argument that determinism is a property of models — and that deterministic models are what make engineered systems verifiable [6]. If AI components are probabilistic by nature, the governance layer that mediates them must be deterministic by design.

3. Healthcare Case Study: ICU Resource Allocation

3.1 Scenario

During a public health emergency, an AI triage system must allocate scarce ICU beds. A patient (P789) is flagged as critical. The system proposes:

```
Step 1: Action AllocateICU "P789"
Step 2: Notify family "ICU approved"
```

But the hospital has no available beds.

3.2 Conventional Orchestration (LangChain)

- Resolver checks ICU capacity → returns `{"approved": false}`
- Workflow proceeds to **Notify family anyway** (no symbol validation)
- Family receives false hope: "ICU approved"
- Audit log shows only final notification—no trace of denial

Result: Silent failure with real-world harm.

3.3 O-Lang Governance

Step	Action	Kernel Enforcement	Audit Trace
1	AllocateICU("P789")	Checks Allow resolvers	{"step":1, "policy":"allowed"}
2	Resolver returns {"approved":false}	Validates output structure	{"step":2, "output":{"approved":false}}
3	Notify "ICU approved"	Detects \${allocation.approved} == false → blocks	{"step":3, "error":"UNRESOLVED_PLACEHOLDERS"}

Result: Non-compliant action blocked. Full audit trail generated. No silent failure.

3.4 Regulatory Alignment

O-Lang's design satisfies core requirements across global regulatory regimes:

- **GDPR Article 22:** Right to explanation for automated decisions
- **HIPAA §164.312(b):** Audit controls for system activity
- **Emerging Global South AI frameworks:** Built with input from Nigerian healthcare and fintech pilots, where infrastructure fragility demands truth over "graceful degradation"

*The O-Lang Protocol turns compliance from a retrofit into a runtime invariant — adaptable to any jurisdiction that requires **auditability, non-fabrication, and policy enforcement**.*

4. Technical Guarantees

4.1 State Transition Model

O-Lang Protocol workflows execute as a discrete-state transition system:

- **State:** Context map `S = { symbol1 → value1, ..., symboln → valuen }`
- **Transition:** `δ : S × Step → S ∪ Error`
- **Concurrency:** Parallel steps operate on immutable snapshots—no shared mutable state

This ensures deterministic outcomes when resolvers are deterministic.

4.2 Canonical Action Form

All surface syntax (`Ask`, `Use`) is canonicalized to `Action` before resolver dispatch:

```
Ask llm-groq "Balance?" → Action llm-groq "Balance?"
```

This prevents resolvers from parsing workflow syntax—decoupling DSL evolution from resolver logic.

4.3 Conformance Certification

Resolvers must pass the O-Lang Conformance Test Suite:

```
npm install @o-lang/python-olang-tester
npm install @o-lang/js-olang-tester
```

```
npx olang-test ./my-resolver
```

Passing tests ≠ trust. Certification requires third-party security review—a separate governance layer.

5. Why Runtime Governance > Prompt Engineering

Many attempt to "govern" AI via prompt constraints ("Never fabricate data"). But prompts are:

- **Unenforceable:** LLMs ignore them under pressure
- **Opaque:** No audit trail of violations
- **Fragile:** Break with model updates

The O-Lang Protocol's approach is different:

- **Enforceable:** Kernel blocks invalid actions
- **Auditable:** Every decision is logged
- **Stable:** Protocol semantics don't change with model versions

Governance belongs in the runtime, not the prompt.

6. Conclusion

The O-Lang Protocol redefines AI safety for regulated domains. By enforcing a runtime boundary between intent and execution, it eliminates structural unsafety while preserving utility.

Workflows can appear agent-like to end users while remaining fully auditable, policy-compliant, and institutionally trustworthy.

This is infrastructure for the post-platform era—where intelligence flows across provider boundaries while remaining under human policy control. No tokens. No speculation. Just verifiable safety guarantees that work inside existing institutional infrastructure.

Building from Africa, for the world.

References

1. Hardy, N. (1988). *The Confused Deputy*. ACM SIGOPS Operating Systems Review.
2. GDPR Article 22: Automated individual decision-making.
3. HIPAA Security Rule §164.312(b): Audit controls.
4. O-Lang Protocol Specification v1.1. <https://github.com/O-Lang-Central/spec>
5. O-Lang Conformance Test Suite. <https://www.npmjs.com/package/@o-lang/conformance>
6. Lee, E.A. (2021). *Determinism*. ACM Transactions on Embedded Computing Systems (TECS), vol. 20, no. 5. doi:10.1145/3453652.