

# Invariant Governance™

## A Structural Architecture for Deterministic Control of Autonomous Systems

**Author:** Joseph Homer Mills **Date:** February 2026 **Version:** 2.0 **Status:** Patent Pending  
(USPTO App. No. 19/533,191; PCT/US26/15432)

---

*This document describes the Invariant Governance architecture, a structural framework for enforcing deterministic safety guarantees over autonomous systems capable of producing irreversible side effects. The architecture is the subject of pending United States and international patent applications and is available as open-source software under the Apache License 2.0.*

---

### Abstract

Autonomous systems — from AI-powered trading agents to robotic surgical assistants to orbital satellites — increasingly operate at speeds and scales where human-in-the-loop approval is physically impossible. Yet these same systems carry the highest consequences for failure. Existing governance approaches rely on behavioral constraints: prompt engineering, reinforcement learning from human feedback, access control lists, and software-defined policy engines. All of these can be bypassed, jailbroken, or silently degraded.

Invariant Governance introduces a fundamentally different approach. Rather than telling autonomous systems what they *should not* do, it architecturally eliminates the *ability* to do it. The system structurally decouples authority (who decides), execution (who acts), and observation (who watches) into three isolated planes connected only by cryptographic artifacts. The result is a governance architecture where safety guarantees are enforced by the absence of capability, not the presence of policy — making them, when correctly implemented, deterministic, auditable, and immune to privilege escalation.

---

## 1. The Problem: Why Existing AI Governance Fails

The AI industry is undergoing a phase transition. For the past decade, AI systems have operated as advisors — suggesting, summarizing, generating. But the current generation of AI agents *acts*. They execute trades, deploy code, send communications, manage infrastructure, and control physical systems. This shift from suggestion to action fundamentally changes the governance problem.

### 1.1 The Latency-Integrity Paradox

The systems that require the highest governance integrity are precisely the systems where human approval introduces unacceptable latency.

Power grid frequency synchronization requires sub-millisecond response. High-frequency trading operates at microsecond granularity. Orbital station-keeping demands continuous real-time adjustment. Autonomous vehicle collision avoidance has zero tolerance for round-trip authorization delays.

Traditional governance models force a choice: operate at machine speed with no governance, or operate with governance at human speed. Neither is acceptable. Invariant Governance resolves this paradox by enabling machine-speed operation within pre-authorized, cryptographically bounded safety envelopes — without requiring per-action human approval.

## 1.2 The Observer Effect

In most existing architectures, monitoring systems possess the ability to trigger automated responses, feed data back into control loops, or modify system state. This creates a dangerous feedback coupling: the monitoring system itself becomes an attack surface.

An adversary who compromises a monitoring system with execution authority gains control over the entire system. A monitoring agent that can “automatically remediate” a detected issue is, by definition, an execution agent wearing a monitoring hat.

Invariant Governance enforces a principle borrowed from physics: measurement must not influence the measured. The observation layer is structurally incapable of modifying the systems it observes — not because of a policy restriction, but because the necessary interfaces simply do not exist.

## 1.3 The Fragmentation Attack (Salami-Slicing)

Consider an AI trading agent authorized to execute transactions up to \$10,000 each. A per-action authorization check will approve every transaction. But what if the agent executes 10,000 transactions of \$9,999 each? Individually, every action is authorized. Collectively, the agent has moved \$99.99 million — potentially far beyond its intended mandate.

This is the fragmentation attack: decomposing a catastrophic action into thousands of individually benign micro-actions that collectively produce an unauthorized outcome. Every per-action governance system is blind to it. Rate limiters catch velocity, not cumulative magnitude. Access control validates identity, not intent drift.

The fragmentation attack is not theoretical. It is the precise mechanism behind historical financial frauds, data exfiltration campaigns, and the predicted failure mode of autonomous agents operating under per-action authorization.

## 1.4 Why Authentication Is Not Enough

Modern security architectures answer the question: “Who are you?” They verify identity through credentials, certificates, and tokens. But an authenticated AI agent with valid credentials can still:

- **Hallucinate** an action that appears consistent with its mandate but produces unintended consequences
- **Drift** from intended behavior through accumulated small deviations over time
- **Be manipulated** through adversarial inputs that exploit the gap between authorization and intent

Authentication establishes identity. Authorization establishes permission. Neither establishes *structural constraint*. An authenticated, authorized agent can still act outside its intended boundaries if the enforcement mechanism relies on the agent’s own compliance.

What is needed is not better authentication or finer-grained authorization, but a governance architecture where unauthorized action is structurally impossible — regardless of the agent’s intent,

credentials, or internal state.

---

## 2. The Architecture: Three Planes, One Invariant

Invariant Governance is built on a single architectural principle: **no single component should possess both the authority to decide and the capability to act**. This principle is enforced through structural decomposition into three isolated planes.

### 2.1 The Authority Plane

The Authority Plane contains the **Governance Kernel** — the sole decision-making component in the architecture. The Kernel evaluates intent proposals from autonomous agents, checks them against policy and cumulative state, and issues cryptographically signed approval receipts when proposals meet all governance requirements.

The critical property of the Governance Kernel is **structural incapacity for execution**. The Kernel possesses no API keys for external services. It has no network routes to execution targets. It holds no credentials that could be used to perform actions. It operates within a hardware-isolated trusted execution environment.

Even if the Governance Kernel is fully compromised — every line of code replaced by an attacker — it *cannot execute* any action in the external world. It can only issue or refuse to issue signed approval receipts. The damage from a compromised Kernel is limited to wrongly approving actions (which must still pass the Execution Gate’s cryptographic signature verification) or wrongly refusing actions (a denial of service, not a safety violation).

### 2.2 The Execution Plane

The Execution Plane contains the **Execution Gate** — a cryptographic enforcement point interposed between every autonomous agent and every execution target. The Gate intercepts all outbound actions, validates them against approval receipts, and releases them only upon successful verification.

The Execution Gate’s critical property is **structural incapacity for decision-making**. It contains no policy logic. It makes no judgment calls. It performs exactly one function: cryptographic verification. It independently computes a hash of the actual execution parameters and compares it to the hash embedded in the approval receipt. Match means proceed. Mismatch means block. No exceptions, no overrides, no fallback.

This separation means that even if an autonomous agent forges or manipulates its own requests, the Execution Gate will block any action whose exact parameters were not pre-approved by the Governance Kernel.

### 2.3 The Observation Plane

The Observation Plane contains the **Advisory Telemetry System** — a structurally isolated monitoring layer that observes system behavior through one-way data channels.

The Advisory Telemetry System’s critical property is **structural incapacity for intervention**. In hardware deployments, physical data diodes (optical one-way devices) prevent any signal from flow-

ing back from the observation layer to the execution or authority layers. In software deployments, read-only log streaming and port mirroring provide the same guarantee.

The Advisory Telemetry System measures system entropy — the divergence between what an autonomous agent was authorized to do and what it is actually doing. It compares observed execution data against a deterministic model of expected behavior. It detects drift, anomalies, and fragmentation patterns. But it can only report. It cannot act.

This structural non-intervention means the observation layer can never be weaponized. Compromising the monitoring system gives an attacker visibility but zero control.

## 2.4 The Sovereign Boundary

The three planes are separated by a **Sovereign Boundary** — a structural, logical, and physical separation that prevents any plane from acquiring the capabilities of another. No shared credentials cross the boundary. No direct communication routes exist between planes except through defined cryptographic interfaces. No component in any plane can escalate its privileges to operate in another.

The Sovereign Boundary is not a firewall. Firewalls are policy-based — they can be reconfigured. The Sovereign Boundary is architectural — the interfaces required for cross-plane operation do not exist and cannot be created without redesigning the system.

---

## 3. Key Mechanisms

### 3.1 Structural Incapacity

The foundational innovation of Invariant Governance is enforcement through *absence* rather than *restriction*.

Traditional security creates a capable system and then restricts it: “This agent has network access but a firewall blocks certain ports.” The agent possesses the capability; the restriction prevents its use. If the restriction fails, the capability is exposed.

Structural incapacity inverts this model. The Governance Kernel does not have network access that is blocked by a firewall. It has *no network interfaces to execution targets at all*. There is nothing to unblock, no policy to bypass, no privilege to escalate. The capability does not exist.

This distinction is critical for autonomous systems. A sufficiently sophisticated AI agent may eventually learn to circumvent policy-based restrictions. It cannot circumvent the absence of a capability.

### 3.2 Cryptographic Approval Receipts

Every authorized action in the system is represented by an **Approval Receipt** — a cryptographically signed, non-fungible authorization artifact. Each receipt contains:

- A cryptographic hash bound to the exact parameters of the approved action
- An expiration timestamp after which the receipt becomes invalid
- A consumption flag ensuring single-use (no replay)
- The entity path identifying who requested the action and under whose authority

- A signature from the Governance Kernel verifiable by the Execution Gate

The Execution Gate independently computes a hash of the actual outbound execution parameters and compares it to the hash in the receipt. Any deviation — even a single changed byte — causes the Gate to block the action. This makes parameter manipulation impossible: the approved action and the executed action must be cryptographically identical.

### 3.3 The Stateful Accumulator (Salami-Slicing Defense)

The **Stateful Accumulator** is a cumulative state tracker maintained by the Governance Kernel. It records the running total of all authorized state changes within a sliding temporal window for each entity in a hierarchical path.

The Governance Kernel enforces a cumulative invariant: the total value of all authorized actions must not exceed a defined safety threshold. This invariant is evaluated on every approval request, using the *historical total* — not just the current request.

This mechanism defeats the fragmentation attack. One thousand requests of \$1 each are treated identically to a single request of \$1,000. The accumulator sees through the decomposition to the cumulative impact. When approaching the safety threshold, the Kernel enters a lockdown state and refuses further approvals regardless of individual request validity.

### 3.4 Dual-Attestation and Mutual Veto

Invariant Governance implements a **mutual veto system** requiring co-signatures from both human operators and autonomous agents.

A human operator signs a strategic work plan — a high-level authorization defining what an autonomous agent is permitted to accomplish. The autonomous agent (Edge Sentinel) signs individual tactical actions — specific executions within the scope of the work plan. Both signatures are required for action. Either party can block by withholding their signature.

This creates bidirectional governance: - The **human can veto the machine** by refusing to authorize a work plan - The **machine can veto the human** by withholding its integrity signature when it detects that executing the requested action would violate governance invariants

This “Bot-in-the-Loop™” capability is as important as “human-in-the-loop.” It prevents the human from being the weakest link — rubber-stamping approvals, succumbing to social engineering, or operating under duress. The system also monitors approval velocity: if a human approves faster than a predefined threshold, the system flags alert fatigue and increases attestation requirements.

### 3.5 Speculative Fluidity Tokens

**Speculative Fluidity Tokens** solve the Latency-Integrity Paradox by enabling pre-authorized autonomous operation within cryptographically bounded safety envelopes.

When the Governance Kernel cannot maintain real-time communication with an Execution Gate — due to network latency, disconnected operations, or throughput requirements — it issues a Fluidity Token containing a risk budget and operational scope. The Execution Gate can then authorize multiple actions locally, decrementing the risk budget with each action, without per-action round-trips to the Kernel.

When the risk budget reaches zero, the Gate blocks all further action until a new token is issued. When connectivity is restored, the Kernel reconciles all actions taken during the disconnected period against the Stateful Accumulator.

This mechanism enables satellite operations during communication blackouts, trading systems during market volatility spikes, and autonomous vehicles during network interruptions — all while maintaining deterministic safety guarantees.

### 3.6 Poison Pill Broadcast

The **Poison Pill Broadcast** is a system-wide emergency halt mechanism. When triggered, a cryptographic signal is transmitted simultaneously to all Execution Gates in an affected jurisdiction. Upon receipt, each Gate immediately destroys its cryptographic key material, enters a lockdown state, and refuses all execution regardless of outstanding approval receipts.

The system remains in lockdown until the Governance Kernel issues new cryptographic keys following human review. This mechanism provides a guaranteed kill switch that cannot be overridden by any autonomous agent, regardless of its credentials or internal state.

### 3.7 Hierarchical Retraction Cascade

Authority in Invariant Governance flows through a **hierarchical entity path** — a tree structure mapping organizational, jurisdictional, or functional scope (for example: /organization/division/team/agent).

A retraction signal issued at any level of the hierarchy immediately invalidates all approval receipts and fluidity tokens at every level below it. A retraction at the organizational level instantly halts every agent in every division, team, and function — with a single signal.

Retraction always supersedes approval. Even a cryptographically valid, unexpired, unconsumed approval receipt is blocked if a retraction signal exists for any ancestor in the entity path. Retractions are recorded in an immutable ledger to prevent rollback or repudiation.

### 3.8 The One-Way Mirror (Advisory Telemetry)

The Advisory Telemetry System operates as a **one-way mirror** — it sees everything but touches nothing. It ingests telemetry data through physically or logically one-way channels, computes a degradation score by comparing observed behavior against expected behavior, and reports anomalies to human operators.

Its outputs are strictly advisory. It cannot trigger automated remediation. It cannot modify governance policy. It cannot issue or revoke approval receipts. This structural constraint ensures that the observation layer can never become an unintended execution channel — a common vulnerability in systems where “monitoring with auto-remediation” creates a hidden backdoor for system-wide control.

### 3.9 The Hardware HAL Gate

For critical physical infrastructure — power grids, water systems, industrial machinery, medical devices — Invariant Governance extends to the hardware level through the **HAL Gate**.

The HAL Gate is a physical enforcement device containing two independent logic modules: a governance module that verifies cryptographic approval receipts, and a physics module that verifies commanded parameters fall within physically safe operating bounds. Only when both modules independently approve does the device activate a galvanically isolated relay connecting digital control signals to physical actuators.

Galvanic isolation provides an electrical air gap between the digital control domain and the physical actuator domain. Even if all software in the system is compromised, the physical actuator cannot be driven outside safe parameters because the hardware logic module independently enforces physics-based constraints.

---

## **4. Application Domains**

Invariant Governance is domain-agnostic by design. The same three-plane architecture applies wherever autonomous systems produce irreversible side effects. The following domains illustrate the breadth of applicability.

### **4.1 Financial Trading and Settlement**

The Execution Gate deploys as a circuit breaker between trading algorithms and exchange matching engines. The Stateful Accumulator tracks cumulative position exposure, preventing flash crashes caused by runaway algorithms. Fluidity Tokens provide risk budgets for high-frequency operation. Shadow telemetry from independent market data feeds detects execution drift.

### **4.2 Clinical Healthcare and Nurse Handoff**

Each nurse handoff interface functions as an Execution Gate. A primary care plan serves as the strategic work plan. Information entropy accumulates with each successive handoff as care plan fidelity naturally degrades. Independent biometric and vitals monitoring (shadow telemetry) provides ground truth, triggering escalation when observed patient state diverges from the care plan.

### **4.3 Orbital and Disconnected Operations**

For orbital satellites, ground-based infrastructure, or any system operating across high-latency or intermittent communication links, Fluidity Tokens enable pre-authorized autonomous operation during disconnected periods. Full governance reconciliation occurs upon connectivity restoration, with the Stateful Accumulator absorbing all actions taken during the disconnected window.

### **4.4 National Infrastructure**

At the national governance scale, the Governance Kernel operates as a sovereignty layer governing critical infrastructure — power grids, water systems, financial systems. The HAL Gate enforces physics-based invariants at the hardware level, preventing silent infrastructure poisoning such as power grid frequency desynchronization or water supply contamination. Hierarchical retraction enables federal-level emergency halt across all subordinate jurisdictions.

#### **4.5 Blockchain and Smart Contract Governance**

The Execution Gate deploys as a smart contract guardrail between a blockchain virtual machine and the public ledger. A private key vault is interlocked with a policy engine: the signing key required to commit a transaction is accessible only after governance verification. Non-compliant transactions are blocked before reaching the immutable ledger — critical because blockchain transactions cannot be reversed after commitment.

#### **4.6 Autonomous Logistics and Swarm Operations**

For fleets of autonomous vehicles, drones, or robots operating in shared physical space, the Governance Kernel issues Fluidity Tokens specifying speed limits, altitude ceilings, and operational boundaries for each agent. A systemic density invariant defines the maximum safe concentration of autonomous agents in any given area, preventing collision, congestion, and airspace conflicts.

#### **4.7 Enterprise Contract Execution**

AI-powered contract execution agents negotiate and execute commercial agreements within governance constraints. The Stateful Accumulator tracks total contractual exposure across a temporal epoch. The Governance Kernel is structurally inhibited from approving any contract that would cause cumulative exposure to exceed a human-ratified threshold — regardless of the individual contract’s validity.

#### **4.8 Government Benefits Issuance**

The hierarchical entity path maps to jurisdictional scopes (federal, state, county, agency). Retraction at a federal or state level immediately invalidates all subordinate payment authorizations. The Advisory Telemetry System monitors for anomalous issuance patterns — sudden approval rate spikes, geographic clustering, temporal correlation with policy changes — triggering advisory alerts for human review.

---

### **5. Why Now**

#### **5.1 The Agent Transition**

Every major AI company is shipping autonomous agents into production. These agents don’t just generate text — they execute code, make API calls, manage cloud infrastructure, send communications, and interact with financial systems. The transition from AI-as-advisor to AI-as-actor is happening now, and it is irreversible.

Yet the governance infrastructure for these agents does not exist. Current approaches rely on behavioral constraints — prompt instructions, RLHF training, content filters — that are fundamentally advisory. They represent the autonomous system’s best effort at compliance, not a structural guarantee.

#### **5.2 The Regulatory Imperative**

The EU AI Act mandates risk management frameworks for high-risk AI systems. The NIST AI Risk Management Framework calls for measurable governance controls. Emerging executive orders

and sector-specific regulations (financial, healthcare, defense) increasingly require demonstrable, auditable governance over autonomous systems.

Invariant Governance provides what these regulations demand: deterministic, cryptographically verifiable proof that autonomous actions were authorized, bounded, and auditable — not probabilistic assurance that an AI model was trained to behave.

### 5.3 The Accountability Gap

When an autonomous agent causes harm, the current accountability model is unclear. Was it the developer’s fault for inadequate training? The deployer’s fault for insufficient guardrails? The operator’s fault for improper oversight?

Invariant Governance creates an unambiguous accountability chain. Every action has a signed approval receipt traceable to a specific policy, a specific human authorization, and a specific governance evaluation. The Stateful Accumulator provides a complete history of cumulative authorization. The immutable retraction ledger provides non-repudiable evidence of governance decisions. This transforms AI accountability from a legal ambiguity into an engineering artifact.

---

## 6. Comparison to Existing Approaches

Capability	Prompt Guardrails	RLHF / Constitutional AI	RBAC / IAM	API Rate Limiting	Circuit Breakers	Invariant Governance
Enforcement type	Behavioral	Behavioral	Policy	Threshold	Threshold	<b>Structural</b>
Bypassable by privileged user	Yes	Yes	Yes	Yes	Yes	<b>No</b>
Defeats jailbreaking	No	Partially	N/A	N/A	N/A	<b>Yes (structural)</b>
Cumulative risk tracking	No	No	No	Rate only	No	<b>Yes (Stateful Accumulator)</b>
Cryptographic audit trail	No	No	Partial	No	No	<b>Yes (Approval Receipts)</b>
Hardware enforcement	No	No	No	No	No	<b>Yes (HAL Gate)</b>
Hierarchical revocation	No	No	Limited	No	No	<b>Yes (Retraction Cascade)</b>

Capability	Prompt Guardrails	RLHF / Constitutional AI	RBAC / IAM	API Rate Limiting	Circuit Breakers	Invariant Governance
Structural non-intervention (observation)	No	No	No	No	No	Yes (One-Way Mirror)
Machine-speed autonomous operation	N/A	N/A	Yes	Degraded	Degraded	Yes (Fluidity Tokens)
Mutual veto (human + machine)	No	No	No	No	No	Yes (Dual-Attestation)

## 7. Licensing and Adoption

### Open Source

The Invariant Governance SDK is released as open-source software under the **Apache License 2.0**. The complete governance framework — all three components, protocol specifications, and reference implementations — is free to use, modify, and distribute. The Apache 2.0 license includes a patent grant for all users of the open-source software, as described in Section 3 of the License.

- **Source:** [github.com/utahbroker/invariant-governance](https://github.com/utahbroker/invariant-governance)
- **Package:** @invariant-governance/core

### Commercial Licensing

Commercial licenses are available for organizations requiring patent coverage beyond the Apache 2.0 grant, dedicated support, certified builds, or compliance documentation for regulated industries.

### Patent Status

The Invariant Governance architecture is protected by pending patent claims:

- U.S. Patent Application No. 19/533,191
- PCT International Application No. PCT/US26/15432
- Related continuation applications

### Contact

For licensing inquiries, partnership opportunities, or technical discussion:

**Email:** [invariant@holladaylabsip.com](mailto:invariant@holladaylabsip.com) **Patent Status:** Pending, United States Patent and Trademark Office

---

*Patent Pending. This document describes the Invariant Governance architecture at a conceptual level for informational purposes. Implementation details are protected under pending patent claims. Use of the open-source SDK under the Apache License 2.0 includes a patent grant as described in Section 3 of the License.*

---

Copyright 2026 Holladay Labs IP, LLC. Licensed under Apache License 2.0.