# Establishing Structural Execution Boundaries for Irreversible AI Actions: The WRS Framework

Jing (Linda) Liu

ORCID: 0009-0002-1681-8563

Independent Researcher

[linda@winston-battery.com]

January 22, 2026

## Abstract

As artificial intelligence systems increasingly transition from decision support to autonomous execution, contemporary AI governance frameworks face a critical structural gap. Existing approaches—ranging from alignment and constitutional constraints to risk scoring and human-in-the-loop oversight—largely assume that execution is permissible once a decision has been produced. This assumption becomes insufficient in systems where execution can trigger irreversible physical, kinetic, or systemic consequences [9], [11].

This paper introduces the World Reliability Ruleset (WRS), a veto-based execution boundary framework designed to govern execution itself rather than decision quality or optimization outcomes. WRS formalizes execution as a binary authorization state governed by a default-block posture: execution is permitted if and only if all non-negotiable constraints are satisfied. Any single violation is sufficient to trigger an absolute veto, making WRS a non-compensatory and non-probabilistic governance mechanism.

Unlike alignment-centric or principle-based models, WRS does not evaluate intent, confidence, or risk gradients. Instead, it defines a deterministic execution boundary that remains independent of the decision engine's intelligence level, preventing increases in reasoning capability from diluting execution safety boundaries. WRS is applicable across advanced AI systems, including AGI and prospective ASI, as well as non-AI automated systems with irreversible consequences.

WRS anchors accountability at the execution boundary by coupling veto events to verifiable binary audit logs, reducing responsibility diffusion and post-hoc rationalization. Positioned as a domain-agnostic execution primitive, WRS complements judgment-layer architectures such as the Linda Energy Reliability Architecture (LERA) while remaining independently deployable. Together, they support a layered governance model in which epistemic exploration may scale without bound, while kinetic manifestation remains strictly constrained by structurally enforced physical permissibility.

This work offers a structural alternative to probabilistic AI safety approaches by introducing a non-compensatory execution-permissibility primitive, providing a governance framework that prioritizes execution authority over decision quality.

## 1.Introduction: The Execution Gap in AI Governance

As artificial intelligence systems advance from analytical tools to agents capable of autonomous action, the focus of AI governance has largely remained anchored in the domain of decision-making. Alignment techniques, constitutional constraints, risk assessment frameworks, and human-in-the-loop (HITL) oversight mechanisms [3], [7], [8] all seek to influence how systems reason, choose, and justify their outputs. Implicit in these approaches is a shared assumption: once a decision has been produced through an acceptable process, execution may proceed as a natural extension of reasoning.

In practice, this "decision-then-execute" assumption is already being stressed by real operational systems that act on the physical world under tight time and coupling constraints. Automated grid dispatch can rapidly redistribute load across interconnected nodes, where a locally "reasonable" adjustment may initiate a cascade that is not recoverable by later correction. Clinical and life-support devices increasingly rely on software-mediated parameter changes, where an incorrect execution is not merely an error but a physiological state transition that may not be reversible within the required time window. Autonomous mobility systems translate perception and planning into kinetic commitment, where inertia turns a late correction into a post-event explanation. And energy systems—battery charging/discharging, thermal control, or protective cutoffs—convert abstract control decisions into stored or released energy that, once uncontained, cannot be rolled back by better reasoning. Across these domains, the governance problem is not only whether the decision was "good," but whether execution should have been possible at all.

This assumption holds in domains where outcomes are reversible, errors are correctable, and system behavior can be iteratively refined. However, it breaks down in environments where execution produces irreversible physical, kinetic, or systemic consequences. In such systems, the transition from decision to execution is not merely procedural; it is a categorical transformation. Once executed, the system irreversibly alters the state of the world, rendering post-hoc correction, re-alignment, or explanation structurally insufficient [10].

This paper argues that contemporary AI governance frameworks systematically under-theorize this transition. While significant effort has been devoted to shaping intelligence—how systems learn, optimize, and align—comparatively little attention has been paid to governing execution

as an independent and non-negotiable phase. This structural omission constitutes what we term the *execution gap* in AI governance.

## 1.1 Decision-Centric Governance and Its Structural Limits

Most prevailing AI governance models are decision-centric by design. Alignment research seeks to ensure that systems internalize human values or objectives. Constitutional and principle-based approaches constrain decision-making through normative rules. Risk management frameworks quantify uncertainty, while human-in-the-loop (HITL) mechanisms reinsert human judgment into critical decision points. Despite their differences, these approaches converge on a common premise: governance operates primarily by shaping or supervising decisions.

In decision-centric models, execution is implicitly treated as derivative. Once a decision satisfies alignment criteria, risk thresholds, or human approval, execution is presumed permissible. Safety, in this view, is achieved by improving the quality of decisions rather than by constraining the act of execution itself.

This premise becomes fragile as systems scale in autonomy, speed, and complexity. As reasoning capability increases, so does the system's capacity to justify its actions internally—often with increasing confidence. Yet confidence, optimization performance, or compliance with abstract principles do not neutralize physical risk. In systems with high consequences, no amount of decision quality can compensate for an unsafe execution event.

Risk-scoring and threshold-based approaches further obscure this limitation. By treating harm as a probabilistic outcome to be minimized, they implicitly assume that sufficiently low risk legitimizes execution. However, irreversibility collapses this logic. Certain outcomes are unacceptable not because they are likely, but because they are possible. This constitutes a **structural blindness** in decision-centric models: they attempt to manage the certainty of physical harm with the uncertainty of probabilistic confidence.

**Table 1. Decision-Centric vs. Execution-Centric Governance**

| Dimension | Decision-Centric Governance | Execution-Centric Governance (WRS) |
|---|---|---|
| Primary focus | Decision quality | Execution permissibility |
| Control mechanism | Alignment, principles, risk thresholds | Structural veto |
| Risk model | Probabilistic, compensatory | Boolean, non-compensatory |
| Treatment of irreversibility | Assumed manageable | Treated as categorical |

| Governance failure mode | Post-hoc justification | Pre-execution prevention |
| --- | --- | --- |

*Table 1 illustrates the categorical distinction between decision-centric and execution-centric governance models.*

## 1.2 Execution as an Event Horizon

Execution in the physical world is a **non-commutative event** with respect to governance. The state reached after execution cannot be restored through post-hoc reasoning, re-alignment, or further optimization. Execution represents the point at which epistemic flexibility—reasoning about alternatives, intentions, and justifications—meets physical finality. Once crossed, no rearrangement of prior decisions can reverse the material consequences that follow.

This non-commutativity exposes a fundamental asymmetry between reasoning and action. Reasoning is epistemic: it explores, simulates, and evaluates possible futures. Execution is kinetic: it commits energy, triggers motion, and manifests safety-critical change in the physical or systemic world. Treating execution as a mere continuation of decision-making conflates these two domains and obscures the moment at which governance must transition from advisory influence to categorical constraint.

From a governance perspective, execution functions as an **event horizon**. It marks the last point of causal intervention at which harm can be structurally prevented. Beyond this boundary, the governance regime necessarily shifts—from prevention to forensics, from control to attribution, and from responsibility avoidance to responsibility assignment. Human oversight, explanation, or accountability invoked after execution may clarify what occurred, but it cannot undo what has been irreversibly enacted.

This distinction reveals why governance mechanisms that operate exclusively upstream—at the level of decision formation—are insufficient in systems with high consequences. Human presence or approval does not constitute control unless it is coupled to a structurally enforced ability to block execution itself. Without such enforcement, governance risks degenerating into symbolic safeguards: visible, reassuring, yet incapable of preventing physical harm.

The execution gap, therefore, is not a failure of intelligence, alignment, or ethical reasoning. It is the absence of a formal boundary that distinguishes what may be reasoned about from what may be enacted. Addressing this gap requires treating execution not as the final step of intelligence, but as the **first step of physical risk**—one that demands its own governance primitives, independent of decision quality, intent, or confidence.

The World Reliability Ruleset (WRS) is a ruleset-based governance framework that enforces non-negotiable execution boundaries through structural constraints. This paper focuses on the structural rationale for treating execution as a veto-eligible state and for separating execution permissibility from decision quality. The companion document—"The World Reliability

Ruleset (WRS): A Technical Specification Supporting the Structural Execution Boundary Framework" [2]—provides the canonical rule text, formal technical standards, trigger semantics, and conformance conditions. Importantly, WRS does not intervene in an algorithm's internal training or reasoning logic; it functions exclusively as a hard gate at the execution boundary. Operational definitions and implementation considerations for execution maturity, emergency override mode switching, and responsibility anchoring are provided in the companion technical specification [2].

## 2. Positioning and Architectural Scope

WRS is intentionally designed as an execution boundary framework rather than a judgment or governance architecture. Its validity does not depend on how judgment is formed, who performs it, or through which institutional or technical process execution is authorized for consideration. The sole prerequisite for WRS applicability is the existence of a candidate execution whose release may produce non-recoverable state consequences.

In this sense, WRS is **architecturally independent** and can, in principle, operate downstream of diverse judgment or authorization frameworks, including legal, organizational, human-led, or hybrid decision structures. This independence is a deliberate design choice, allowing WRS to function as a **domain-agnostic execution boundary** across heterogeneous governance environments, without prescribing how judgment ought to be produced.

Architectural independence, however, does not imply governance completeness. WRS does not evaluate judgment sufficiency, correctness, or legitimacy, nor does it define responsibility formation or authority allocation. It assumes that some form of judgment has already occurred and restricts itself strictly to the execution layer, where irreversible consequences are about to be released. As such, WRS cannot, by itself, constitute a closed governance system.

Within the present work, WRS is positioned to operate downstream of **LERA (Judgment–Governance Architecture)**, which provides a structural mechanism for judgment formation, governance gating, and responsibility anchoring prior to execution consideration. In this configuration, **LERA filters for responsibility, while WRS filters for physical safety**. The two frameworks are complementary but non-overlapping: LERA establishes execution eligibility in principle, whereas WRS constrains execution permissibility in practice.

This separation is not merely conceptual but structural. WRS must not be interpreted as a substitute for judgment architectures, nor as a mechanism for decision-making or authorization. Conversely, judgment architectures—including LERA—do not eliminate the need for execution-layer vetoes once irreversible consequences are at stake. Together, they define distinct but interoperable layers of governance, preserving accountability at both the judgment boundary and the execution boundary.

# 3. Why Existing AI Governance Frameworks Fail at the Execution Boundary

## 3.1 Alignment-Centric Governance: Optimization Without Execution Constraints

A substantial portion of contemporary AI governance research is centered on alignment. From value learning and reward modeling to preference optimization and constitutional constraints, alignment-centric approaches aim to ensure that intelligent systems pursue objectives consistent with human intent, ethical norms, or institutional goals. These efforts address a critical question: *what should the system optimize for?*

However, alignment-based governance largely operates upstream of execution and does not impose constraints on whether execution should occur at all. Once an aligned decision has been produced—whether through optimized reward functions, policy constraints, or post-training safeguards—execution is typically treated as a permissible continuation of intelligence rather than as a governed state requiring independent validation.

This structure reveals a fundamental limitation. Alignment governs **objective consistency**, not **execution permissibility**. An aligned system may correctly optimize for its intended goals and still produce actions whose execution results in irreversible harm. In such cases, the quality of alignment does not mitigate the consequences of execution; it merely explains why the action occurred.

The distinction becomes clearer when execution is understood as a **categorical transition** rather than a procedural step. **Execution is a non-commutative event in the physical world; it creates a state that cannot be reverted through re-alignment or further training.** Once an action is executed, subsequent improvements in model behavior, updated objectives, or refined preferences cannot undo the resulting physical or systemic outcome.

As systems increase in confidence, autonomy, and operational scope, this limitation intensifies. Alignment mechanisms often enhance a system's ability to justify its actions, produce coherent explanations, and satisfy predefined objectives. Yet this very capacity for justification can reinforce over-permission, encouraging execution precisely because an action appears well-reasoned, aligned, or statistically supported.

In this sense, alignment fosters **justification**, whereas execution governance demands **restraint**. Alignment relies on the system's internal confidence and optimization success; execution governance must honor the constraints imposed by external physical reality. Without explicit execution boundaries, alignment functions as an explanatory layer rather than a limiting one. Even when aligned objectives are specified, concrete safety failure modes - such as reward hacking, unsafe exploration, and robustness breakdowns - remain well-documented in high-autonomy systems [4].

Moreover, alignment frameworks are inherently grounded in continuous, gradient-based improvement. Such approaches excel at managing trade-offs and preferences but struggle to represent the discontinuities introduced by high consequences. Where execution outcomes cannot be undone, compensated, or meaningfully corrected, incremental improvements in alignment offer no structural safeguard.

As a result, alignment-centric governance collapses the distinction between deciding and acting. Once a decision satisfies alignment criteria, execution proceeds by default. This conflation obscures the fact that execution marks a transition from abstract reasoning to irreversible consequence—a transition that cannot be governed by optimization quality alone.

In the absence of explicit execution constraints, alignment becomes a mechanism of rationalization rather than protection. It explains *why* an action appears reasonable, but does not determine *whether* that action must be categorically blocked when irreversible consequences are at stake. This structural gap persists regardless of how advanced, accurate, or well-aligned the system becomes.

## 3.2 Human-in-the-Loop Oversight: Supervision Without Structural Authority

Human-in-the-loop (HITL) oversight is widely regarded as a cornerstone of responsible AI governance. By inserting human review, approval, or intervention into automated decision pipelines, HITL frameworks aim to preserve human involvement over consequential actions while benefiting from machine efficiency and scale. In principle, the presence of a human decision-maker is assumed to mitigate the risks associated with autonomous execution.

In practice, however, HITL mechanisms frequently operate as **supervisory layers rather than execution authorities**. Human review is often advisory, confirmatory, or procedurally constrained, embedded within workflows optimized for continuity and throughput. As a result, the human role tends to validate decisions rather than to structurally constrain execution.

This distinction is critical. **Human presence is a symbolic rather than a functional constraint unless the veto is formalised as a blocking primitive in the execution logic.** Without such a primitive, human approval does not interrupt execution pathways; it merely annotates them. Authority exists only insofar as the system can be categorically prevented from acting, independent of confidence, urgency, or procedural momentum.

Temporal misalignment further undermines the effectiveness of HITL governance. Human intervention frequently occurs after key execution commitments have already been initiated—resources allocated, processes activated, or external systems engaged. At this stage, withholding execution imposes escalating operational, financial, or institutional costs, implicitly biasing reviewers toward approval. The execution path remains open, while the human role becomes increasingly constrained.

In systems involving irreversible consequences, this misalignment becomes structural rather than incidental. **HITL oversight often functions as a post-mortem authority acting within a nominal pre-execution window**, where review occurs too late to meaningfully alter the outcome. Once execution readiness has been established, human intervention serves to attribute responsibility rather than to prevent release.

Cognitive and organizational dynamics further weaken HITL effectiveness. Repeated exposure to system recommendations can induce automation bias, particularly when systems exhibit high apparent accuracy or internal confidence. Over time, human operators are positioned as exception handlers rather than as execution gatekeepers, reinforcing the assumption that action should proceed unless an anomaly is detected.

As with alignment-centric governance, HITL frameworks tend to conflate decision validation with execution permissibility. Once a human has reviewed or endorsed a decision, execution is treated as legitimate by default. This conflation obscures the fact that human approval, absent structural veto authority, does not resolve the risks associated with irreversible execution.

In the absence of enforceable execution boundaries, HITL oversight functions primarily as a mechanism of accountability attribution rather than prevention. It may clarify who reviewed or approved an action after the fact, but it does not reliably block execution when irreversible consequences are at stake. As a result, human-in-the-loop governance, while valuable for transparency and responsibility tracing, remains structurally insufficient as a safeguard against execution-level failure.

## 3.3 Risk Scoring and Probabilistic Governance: Gradients Where Discontinuities Exist

Risk-based governance occupies a central position in contemporary AI safety and regulatory practice. Through probabilistic modeling, impact assessments, confidence thresholds, and escalation rules, such approaches aim to quantify uncertainty and manage harm by balancing likelihood against severity. In many domains, this framework provides a pragmatic means of prioritizing attention and allocating safeguards.

However, probabilistic governance presupposes a **continuous risk landscape**, in which outcomes can be compared, trade-offs evaluated, and marginal improvements meaningfully reduce harm. This assumption breaks down in systems where execution produces **irreversible consequences**. In such contexts, risk does not scale smoothly with probability, and the distinction between acceptable and unacceptable outcomes cannot be captured by gradients or thresholds.

Irreversibility introduces a structural discontinuity. Once an execution event occurs, the resulting state cannot be undone, compensated, or meaningfully corrected through post hoc intervention. **Irreversible harm is not a peak in a risk landscape; it is the collapse of the**

**landscape itself.** The expected-value logic underpinning probabilistic risk management—where low-probability events may be tolerated in exchange for aggregate benefit—fails to represent this collapse. A small probability multiplied by an irreversible outcome does not yield a "small risk"; it yields a permanent state change.

Despite this, risk-scoring frameworks routinely authorize execution based on numerical thresholds or confidence intervals. Actions are permitted not because they are categorically safe, but because they fall below a predefined tolerance. This transforms execution into a statistical entitlement: as long as modeled risk remains within acceptable bounds, execution proceeds by default. Governance thus becomes a question of calibration rather than prohibition.

This framing obscures a critical asymmetry. Risk models operate in **epistemic space**, representing beliefs, predictions, and uncertainty, while execution operates in **physical space**, where outcomes are concrete and irreversible. No increase in predictive accuracy or model calibration can bridge this gap. A perfectly estimated probability does not mitigate the consequence of an irreversible execution; it merely anticipates it.

Furthermore, probabilistic governance normalizes rare but catastrophic outcomes through repetition. As systems are deployed at scale, low-probability failures accumulate across executions, transforming statistical exceptions into eventual certainties. Yet threshold-based governance treats each execution independently, ignoring the inevitability of collapse in systems lacking absolute constraints.

At the execution boundary, this reveals a categorical mismatch. **Risk thresholds manage uncertainty, while execution governance must govern possibility.** The former asks *how likely* an outcome is; the latter must ask *whether the outcome should be possible at all*. When irreversible consequences are at stake, permissibility cannot be derived from probability alone.

As a result, probabilistic governance frameworks systematically underperform at the execution boundary. They offer refinement without restraint, confidence without containment, and optimization without veto. In the absence of mechanisms capable of categorically blocking execution, risk scoring functions as a justification layer rather than a protective one, leaving irreversible outcomes to be addressed only after they have already occurred.

## 3.4 Compliance-Driven Governance: Documentation Without Execution Boundaries

Compliance- and principle-based governance frameworks play a central role in contemporary AI regulation and organizational risk management. Through ethical guidelines, safety principles, documentation requirements, audits, and certification processes, these approaches aim to ensure that automated systems adhere to legal obligations, institutional standards, and societal expectations. Their primary contribution lies in establishing norms of conduct and mechanisms for accountability.

However, compliance frameworks predominantly operate at the level of **eligibility**, not **permissibility**. They govern whether a system, organization, or deployment context meets the criteria to operate, but they do not govern whether a specific execution event should be allowed to occur. In other words, compliance frameworks manage eligibility for deployment, but they do not regulate the permissibility of individual execution events at the moment when high consequences may be released.

This distinction exposes a structural limitation. Compliance mechanisms are effective at answering questions such as *who was responsible*, *which requirement was satisfied*, or *whether due process was followed*. They are far less effective at addressing a more fundamental governance question: *should this execution have been possible at all?* When irreversible harm occurs, post hoc accountability does not mitigate the outcome; it merely documents the conditions under which it happened.

Principle-based governance further reinforces this gap. High-level principles—such as transparency, fairness, safety, or human oversight—provide valuable normative orientation but lack operational specificity at the execution boundary. Principles articulate desirable properties of systems and processes, not categorical constraints on action. As a result, they are frequently interpreted flexibly, balanced against competing objectives, or satisfied through documentation rather than enforcement.

In practice, compliance and principle-based systems tend to legitimize execution through procedural completion. Once audits are passed, certifications obtained, or checklists satisfied, execution is treated as authorized by default. Governance thus becomes permissive: adherence to process substitutes for structural restraint. The absence of explicit blocking conditions allows execution pathways to remain open even in the presence of physically consequential risk.

This limitation becomes acute in systems where execution produces irreversible physical, systemic, or life-critical consequences. In such contexts, **compliance without a veto mechanism is merely the documentation of an inevitable failure**. Responsibility may be clarified after the fact, but the harm itself cannot be undone. Governance operates retrospectively, while execution proceeds unimpeded.

Moreover, compliance regimes typically lack real-time operational authority. Audits, reviews, and certifications are periodic and detached from live execution logic. When execution unfolds rapidly or autonomously, compliance artifacts remain static, offering no mechanism to interrupt or block action at the critical moment.

As a result, compliance and principle-based governance frameworks function primarily as instruments of legitimacy rather than prevention. They establish that systems were built and deployed according to accepted standards, but they do not ensure that execution will be blocked when structural risk materializes. This is not a failure of regulation or ethics, but a mismatch of governance layer.

Preventing high-consequence harm requires governance primitives that operate directly on **execution permissibility**, not merely on deployment eligibility or procedural conformity. Without an explicit execution boundary, compliance assures accountability while leaving execution fundamentally unconstrained.

## 3.5 The Structural Absence of an Execution Boundary

The limitations identified across alignment-centric governance, human-in-the-loop oversight, probabilistic risk management, and compliance-based regulation converge on a common structural deficiency. In each case, governance mechanisms operate either upstream of execution or retrospectively after it has occurred. None impose categorical constraints on execution itself at the moment of **physical commitment**, where abstract reasoning translates into irreversible kinetic or systemic change.

This absence is not incidental. Contemporary AI governance frameworks are predominantly designed to govern **decision quality**, **oversight processes**, **risk estimation**, or **procedural legitimacy**. Execution is treated as an implicit entitlement that follows once these conditions are satisfied. As a result, execution permissibility is inferred indirectly—from confidence, approval, probability, or compliance—rather than governed as a primary and independent condition.

The consequence is a systematic over-permission of action. When execution pathways remain structurally open, governance functions as justification rather than containment. Decisions may be aligned, reviewed, assessed, and documented, yet execution proceeds by default even in scenarios where irreversible harm may result from the release of physical energy or systemic commitment.

What is missing across these frameworks is an explicit **execution boundary**: a governance layer in which execution is treated as a conditional, **veto-eligible state** rather than an assumed continuation of intelligence or process. Crucially, the governance primitives required to enforce such a boundary must be **orthogonal to the system's intellectual objectives**. They must remain invariant regardless of the complexity, confidence, or intent of the decision engine, ensuring that safety constraints cannot be diluted, optimized away, or overridden by advances in intelligence.

Without such orthogonal primitives, no combination of alignment, supervision, estimation, or compliance can reliably prevent execution-level failure. Governance remains bound to the logic of decision-making, while execution—where irreversible consequences materialize—escapes direct constraint.

This structural absence motivates the need for governance mechanisms that operate directly on **execution permissibility** at the moment of physical commitment. The following section introduces a framework designed to fill this gap by formalizing execution vetoes as a foundational governance primitive.

# 4. The World Reliability Ruleset (WRS)

The preceding sections have demonstrated that prevailing approaches to AI governance—whether grounded in alignment, human oversight, probabilistic risk management, or compliance—systematically fail to constrain execution at the point where irreversible consequences materialize. These failures do not arise from insufficient optimization, supervision, estimation, or regulation, but from the absence of a governance layer that treats execution itself as a first-class condition.

This section introduces the **World Reliability Ruleset (WRS)** as a response to this structural gap. WRS is not a decision-making framework, an optimization strategy, or a risk management methodology. It is a **normative execution boundary**: a ruleset designed to govern whether execution may occur at all, independent of how decisions are generated, justified, or approved.

WRS redefines execution **not as the final step of intelligence, but as the first step of physical risk**. At the moment of execution, abstract reasoning is translated into physical commitment, where irreversible kinetic or systemic change may be released. It is at this boundary—not within the reasoning process itself—that governance must become absolute.

At its core, WRS formalizes execution as a **veto-eligible state** rather than an assumed entitlement. Execution is permitted only when all governing constraints are satisfied; any single violation is sufficient to block action categorically. This default-block posture reflects the recognition that once physical commitment occurs, subsequent explanation, optimization, or correction cannot reverse the outcome.

Crucially, WRS operates **orthogonally to intellectual objectives**. Its constraints are invariant with respect to system intelligence, confidence, optimization success, or intent alignment. Whether decisions are produced by narrow models, general intelligence, or superintelligent systems, the execution boundary enforced by WRS remains unchanged. Intelligence may improve reasoning, but it cannot override the conditions under which action is permitted.

The execution boundary defined by WRS constitutes a **non-negotiable interface** between judgment and action. This interface serves as a **formal contract between the reasoning engine and the physical environment**, specifying the conditions under which abstract decisions may be converted into real-world effects. It does not evaluate decision quality, balance trade-offs, or assign responsibility; it governs permissibility at the point of physical commitment.

WRS is therefore agnostic to architecture, training paradigm, and upstream governance regime. It does not compete with alignment frameworks, oversight mechanisms, risk models, or compliance processes. Instead, it constrains them, ensuring that execution remains structurally blocked whenever irreversible consequences are at stake.The remainder of this section specifies the axiomatic foundations, design principles, and normative structure of WRS, establishing execution vetoes as a governance primitive rather than a discretionary safeguard.

## 4.1 Axiomatic Foundations of WRS

The World Reliability Ruleset (WRS) is grounded in a small set of axioms that define non-negotiable constraints on execution. These axioms do not describe how decisions are made, optimized, or justified. Instead, they specify the conditions under which execution may or may not occur once a candidate action reaches the execution boundary.

The axioms are intentionally minimal, discrete, and invariant with respect to system intelligence. Together, they establish execution vetoes as a governance primitive rather than a discretionary safeguard.

---

**Axiom 1：Execution Irreversibility Axiom**

**If an execution may produce irreversible physical, systemic, or life-critical consequences, execution legitimacy cannot be inferred from decision quality, confidence, or optimization success.**

This axiom establishes a fundamental separation between reasoning and consequence. Once execution results in irreversible change, no improvement in alignment, prediction, or justification can retroactively legitimize the action. Execution legitimacy must therefore be evaluated independently of the reasoning process that produced it.

---

**Axiom 2：Default-Block Axiom**

**In the presence of potentially irreversible consequences, execution is structurally blocked by default and may proceed only if no governing constraint is violated.**

This axiom defines the default posture of WRS. Execution is not presumed permissible. Permission must be earned through the absence of violations, not granted through confidence, approval, or optimization. The default-block condition ensures that execution remains constrained unless explicitly cleared.

---

**Axiom 3：Single-Veto Sufficiency Axiom**

**A single violation of any execution constraint is sufficient to invalidate execution permissibility.**

This axiom formalizes the veto principle. Execution permissibility does not arise from aggregate scoring, weighted trade-offs, or threshold satisfaction. Any individual constraint violation is sufficient to block execution categorically. No compensatory reasoning or probabilistic offset is admissible.

This establishes WRS as a non-compensatory framework, in which favorable attributes or high-confidence assessments cannot neutralize a single disqualifying condition.

**Axiom 4：Judgment–Execution Separation Axiom**

**The legitimacy of execution is distinct from, and irreducible to, the legitimacy of judgment.**

This axiom defines the interface between WRS and judgment-based governance architectures. Judgment may determine whether an action is proposed, reviewed, or endorsed, but it does not determine whether execution is permitted. Execution legitimacy requires an independent evaluation at the execution boundary.

**Axiom 5：Orthogonality to Intelligence Axiom**

**Execution constraints enforced by WRS are orthogonal to the system's intellectual objectives and remain invariant with respect to intelligence, confidence, intent alignment, or optimization capability.**

This axiom ensures that execution safety boundaries cannot be optimized away, diluted, or overridden by increases in system intelligence. Whether decisions are produced by narrow automation, artificial general intelligence, or superintelligent systems, the execution constraints imposed by WRS remain unchanged.

**Axiom 6：Physical Commitment Axiom**

**Execution is defined as the point of physical commitment, where abstract reasoning is converted into irreversible kinetic or systemic change.**

Execution constitutes the event horizon of governance, beyond which prevention is structurally impossible. Once this boundary is crossed, no further reasoning, oversight, or corrective intervention can reverse the material consequences. Governance must therefore operate decisively at this boundary, not beyond it.

Together, these axioms establish WRS as a ruleset that governs **execution permissibility**, not decision quality or procedural legitimacy. They define execution as a veto-eligible state and formalize default-blocking as the only governance posture consistent with irreversible consequences.

## 4.2 Design Principles

The axioms defined in Section 4.1 establish non-negotiable constraints on execution. From these constraints follow a set of design principles that characterize how an execution governance system must be structured if it is to remain consistent with WRS. These principles do not prescribe implementation details; they specify invariant properties that any WRS-compliant system must exhibit.

**Principle 1：Execution over Decision**

Execution governance must take precedence over decision evaluation.

WRS-compliant systems do not seek to improve the quality, alignment, or justification of decisions as a prerequisite for safety. Instead, they govern the permissibility of execution independently of how decisions are produced.

This principle follows directly from the Judgment–Execution Separation Axiom. Decision quality may influence what actions are proposed, but it cannot determine whether execution is allowed. Execution is treated as a distinct governance domain, not as an extension of reasoning. **This ensures that the execution boundary remains computationally orthogonal to the decision engine's objective function**, preventing safety constraints from being absorbed, optimized, or traded off within the decision process itself.

## Principle 2：Boolean over Probabilistic Control

Execution permissibility must be represented as a discrete Boolean state rather than a probabilistic or graded measure.

In WRS, execution is either **permitted** or **blocked**. There are no confidence thresholds, risk scores, or acceptable loss margins that can conditionally authorize action. This principle enforces the non-compensatory nature of execution governance and prevents risk normalization through aggregation or averaging.

Probabilistic reasoning may inform judgment, but it cannot arbitrate execution.

## Principle 3：Veto over Optimization

Execution safety must be enforced through veto mechanisms rather than optimization objectives.

WRS does not assign rewards, penalties, or trade-offs to safety conditions. Any attempt to encode execution constraints as optimization terms risks dilution under competing objectives. Instead, WRS treats execution constraints as categorical vetoes: violations block execution without negotiation or compensation.

This principle ensures that safety boundaries remain invariant regardless of optimization pressure or performance incentives.

## Principle 4：Invariance under Intelligence Scaling

Execution boundaries must remain invariant as system intelligence scales.

A WRS-compliant execution boundary does not adapt, relax, or recalibrate in response to increased model capability, confidence, or autonomy. Improvements in reasoning, prediction, or planning do not modify the conditions under which execution is permitted.

This principle operationalizes the Orthogonality to Intelligence Axiom and ensures that advances in intelligence do not erode execution safety.

**Principle 5：Boundary before Authority**

Execution boundaries must be evaluated prior to, and independently from, authority, approval, or responsibility assignment.

Human approval, institutional authorization, or procedural compliance cannot substitute for execution boundary validation. Authority may exist upstream of execution, but it does not override veto conditions at the execution boundary.

This principle prevents execution from being legitimized solely through endorsement or procedural completion.

**Principle 6：Prevention over Remediation**

Execution governance must prioritize prevention over remediation.

WRS assumes that once physical commitment occurs, remediation is structurally limited or impossible. Therefore, governance effort is concentrated at the execution boundary, where blocking remains feasible. Post hoc correction, explanation, or compensation is treated as insufficient for irreversible harm.

This principle aligns governance temporality with the point at which meaningful control can still be exercised.

**Principle 7：Minimalism and Non-Extensibility of Core Constraints**

The core execution constraints of WRS must remain minimal and non-extensible.

Adding complexity to execution rules increases ambiguity and undermines enforceability. Domain-specific considerations may be addressed through constrained extensions, but the core execution principles must remain fixed, discrete, and non-negotiable.

**Any expansion of the core constraints must be treated as a version-breaking modification rather than an interpretation**, preserving the integrity, auditability, and authority of the ruleset.

Together, these principles describe an execution governance regime that is discrete, invariant, and resistant to optimization pressure. They translate the axiomatic foundations of WRS into structural requirements that clearly distinguish execution boundaries from decision processes, compliance procedures, and risk management systems.

## 4.3 Formal Semantics: Execution as a Boolean State

To operationalize the axioms and design principles of WRS, execution must be represented using a formal semantic model that is discrete, non-compensatory, and invariant under optimization pressure. This section defines the semantic status of execution within WRS and explains why execution permissibility cannot be expressed as a continuous, probabilistic, or threshold-based quantity.

### 4.3.1 Execution as a Binary Governance State

Within WRS, execution is modeled as a Boolean governance state:

$$Execution \in \{Permitted, Blocked\}$$

No intermediate, partial, or probabilistic execution states are admissible. An action either crosses the execution boundary or it does not. This binary formulation reflects the structural reality that execution constitutes a discrete transition from abstract reasoning to physical commitment.

Once execution is permitted, the system enters a new physical or systemic state that cannot be reverted through further reasoning, alignment, or optimization. Conversely, when execution is blocked, no amount of confidence, approval, or predicted benefit can partially authorize action.

### 4.3.2 Rejection of Gradient-Based Safety Models

Many contemporary governance and safety frameworks rely on gradient-based representations of risk, confidence, or acceptability. These models assume that safety can be increased incrementally and that sufficiently low risk may justify execution.

WRS explicitly rejects this assumption at the execution boundary. In the presence of irreversible consequences, gradients collapse. Irreversible harm is not a peak in a risk landscape; it is the collapse of the landscape itself. Once physical commitment occurs, marginal improvements in prediction or alignment no longer alter the outcome.

By representing execution as a Boolean state, WRS removes the possibility of negotiating safety through probability, confidence intervals, or expected-value calculations.

### 4.3.3 Non-Compensatory Semantics and Veto Logic

The Boolean execution model directly encodes the non-compensatory nature of WRS. No favorable attribute—such as high expected utility, strong alignment confidence, regulatory approval, or human endorsement—can compensate for a single violated execution constraint.

Formally, let

$$C = \{c_1, c_2, \ldots c_n\}$$

be the set of execution constraints. Execution is permitted if and only if:

$$\bigwedge_{i=1}^{n} eval(c_i) = True$$

If any constraint evaluates to false, execution is categorically blocked. No aggregation, weighting, or thresholding is admissible. This veto logic prevents the normalization of risk through accumulation of favorable signals and ensures that execution remains veto-eligible by semantic construction rather than policy preference.

### 4.3.4 Execution as an Event Horizon

Execution in the physical world is a non-commutative event. The order of operations matters: once execution occurs, subsequent reasoning, explanation, or corrective action cannot restore the prior state.

Execution therefore constitutes the **event horizon of governance**. Beyond this point, prevention is structurally impossible and information feedback loops are irreversibly severed. While computational processes may allow revision or rollback, physical commitment introduces non-recoverable state transitions that no post hoc intervention can negate.

The Boolean semantics of WRS reflect this reality. Execution cannot be partially undone, softened, or compensated after the boundary is crossed.

### 4.3.5 Semantic Separation from Decision and Compliance Layers

The Boolean execution state defined by WRS is semantically distinct from decision confidence, approval status, or compliance certification. These attributes may inform judgment or deployment eligibility, but they do not directly influence execution permissibility.

A decision may be well-justified, well-reviewed, and fully compliant, yet still be blocked at the execution boundary. Conversely, no decision—regardless of endorsement or authority—can bypass a blocked execution state. This establishes a semantic firewall between reasoning, compliance, and action.

### 4.3.6 Auditability and Proof of Compliance

The Boolean execution model enables precise and objective auditability. Because execution resolves to a discrete state—permitted or blocked—the audit trail becomes a **deterministic log of veto triggers**, rather than a subjective interpretation of risk-score fluctuations or confidence estimates.

This structure simplifies proof of compliance. Instead of reconstructing probabilistic reasoning or post hoc justifications, auditors can verify execution legitimacy by inspecting whether any execution constraint was violated at the moment of physical commitment.

Boolean semantics thus preserve governance intent over time, prevent interpretive drift, and provide a clear evidentiary basis for accountability without relying on probabilistic explanations.

By formalizing execution as a Boolean governance state, WRS translates its axiomatic commitments into a semantic structure that is resistant to ambiguity, optimization pressure, and reinterpretation. Execution is governed not by how confident, compliant, or intelligent a system may be, but by whether execution is categorically permissible at the point of physical commitment.

### 4.4 Normative Structure of WRS

The World Reliability Ruleset (WRS) is defined as a normative execution governance system composed of a stable core and constrained domain-specific extensions. This structure preserves the invariance, non-compensatory semantics, and auditability established in the preceding sections, while enabling applicability across heterogeneous physical and automated domains.

### 4.4.1 Core–Domain Separation

WRS is structured around two distinct layers:

- **WRS-C (Core Constraints)**

- **WRS-D (Domain Constraints)**

This separation is foundational. The core constraints define universal execution conditions that apply to all systems subject to WRS, regardless of application domain, system architecture, or level of intelligence. Domain constraints specialize execution governance for particular classes of physical interaction without modifying, weakening, or reinterpreting the core.

All execution constraints—core or domain—are evaluated under the same Boolean, non-compensatory semantics defined in Section 4.3.

### 4.4.2 WRS-C: Core Execution Constraints

WRS-C consists of a minimal, closed set of execution constraints governing execution permissibility in the presence of irreversible consequences. These constraints are:

- **Universal**: applicable across all domains where execution may result in irreversible physical, systemic, or life-critical outcomes.

- **Invariant**: unchanged by system intelligence, optimization capability, or governance context.

- **Non-negotiable**: not subject to trade-offs, weighting, or reinterpretation.

WRS-C does not encode domain knowledge, operational heuristics, or engineering thresholds. Instead, it defines absolute execution conditions that must be satisfied before any domain-specific evaluation may proceed.

If any WRS-C constraint is violated, execution is categorically blocked, independent of downstream considerations.

### 4.4.3 WRS-D: Domain-Specific Constraint Sets

**WRS-D defines a family of constrained execution domains, each representing a domain-specific subset of WRS-D under a shared execution semantics.**
WRS-D defines constrained extensions that tailor execution governance to specific categories of

physical interaction.Each WRS-D module introduces additional execution constraints relevant to a particular domain while inheriting all WRS-C constraints without modification.

Representative domain constraint sets include:

- **WRSE (Energy)**: governing execution involving energy storage, release, conversion, or transmission.

  *For instance, WRSE may specify a non-negotiable veto condition when cell-level temperature gradients exceed safety parameters, independent of external energy demand or system-level optimization objectives.*

- **WRSM (Mobility)**: governing execution involving motion, transport, or kinetic interaction.

- **WRSG (Grid / Infrastructure)**: governing execution involving critical infrastructure and systemic interdependence.

These domain constraint sets are representative rather than exhaustive, and serve to illustrate the structure of WRS-D rather than to enumerate all possible domains.

Domain constraints refine execution permissibility by introducing additional veto conditions. They do not alter execution semantics, default-block posture, or veto logic.

### 4.4.4 Constraint Evaluation and Veto Propagation

Execution permissibility under WRS is determined by the conjunction of all applicable constraints:

$$\textit{Execution Permitted} \Longleftrightarrow \bigwedge_{c \in WRS-C \cup WRS-D} eval(c) = \text{True}$$

This structure ensures that:

- Core constraint violations always propagate vetoes upward.

- Domain constraint violations block execution without exception.

- No constraint—core or domain—can be overridden by authority, confidence, or optimization success.

Veto sufficiency is guaranteed by non-compensatory semantics; evaluation order is logically irrelevant.

### 4.4.5 Domain Extensibility as a One-Way Gate

WRS supports domain extensibility without compromising rule integrity. New WRS-D modules may be introduced to address emerging domains of physical interaction, provided that they:

1. Do not modify, reinterpret, or weaken any WRS-C constraint.

2. Preserve Boolean evaluation semantics.

3. Introduce only additional veto conditions, never permissions.

Structurally, WRS-C functions as a **one-way gate**. Domain constraints operate behind this gate as finer-grained barriers: they may block additional execution paths, but they can never allow an execution that the core has already prohibited.

Any modification to WRS-C or alteration of execution semantics constitutes a **version-breaking change**, not an interpretation or extension. Such changes require formal revision of the ruleset itself.

### 4.4.6 Normative Authority and Compliance Interface

The normative authority of WRS derives from its role as an execution boundary rather than a decision framework. Compliance with WRS is demonstrated through verifiable adherence to execution veto conditions, not through predictive accuracy or risk minimization.

Because execution outcomes resolve to discrete states, WRS enables a clear compliance interface:

- Execution requests are evaluated against WRS-C and applicable WRS-D constraints.

- Veto triggers are logged as discrete, auditable events.

- Permissibility is established through absence of violations, not probabilistic justification.

### 4.4.7 Structural Position within LERA-Compatible Architectures

Within architectures that incorporate judgment-based governance layers—such as LERA [1]—WRS occupies a distinct and non-overlapping role. Judgment determines whether execution should be considered; WRS determines whether execution may occur.

LERA governs **deliberation and responsibility anchoring**.
WRS governs **physical execution and irreversible commitment**.

This separation ensures modularity while enforcing an absolute execution boundary.

### 4.4.8 Summary

The normative structure of WRS formalizes execution governance as a layered, veto-based system composed of immutable core constraints and constrained domain extensions. By separating universality from specialization and extensibility from mutability, WRS preserves governance invariance while enabling broad applicability.

WRS thus functions as a stable execution boundary protocol rather than a domain-specific safety mechanism.

## 4.5 Execution Boundary and Accountability

A central concern in the governance of automated systems is accountability: who bears responsibility when systems act, fail, or cause harm. In conventional governance models, accountability is often entangled with decision-making authority, system intelligence, or human oversight. WRS deliberately disentangles these concepts by anchoring accountability at the execution boundary rather than within the decision process itself.

### 4.5.1 Accountability as a Consequence of Execution, Not Decision

Under WRS, accountability is triggered by **execution**, not by reasoning, recommendation, or approval. Decisions—whether generated by humans, AI systems, or hybrid processes—do not themselves produce irreversible consequences. Execution does.

A system may generate unsafe, incorrect, or even extreme decisions without causing harm, provided those decisions are blocked at the execution boundary. Conversely, a decision that is widely endorsed, well-reasoned, or procedurally compliant may still produce irreversible harm if executed. Accountability therefore arises only when execution is permitted and physical commitment occurs.

This distinction establishes a protected speculative space for exploration, simulation, and judgment, while preserving strict responsibility for real-world effects.

### 4.5.2 The Execution Boundary as a Point of Responsibility Transfer

The execution boundary defined by WRS marks the point at which abstract intent is converted into real-world effect. It is at this boundary that responsibility is transferred from deliberative processes to accountable actors.

When execution is **blocked**, no physical commitment occurs and no responsibility for real-world harm is incurred. When execution is **permitted**, responsibility becomes anchored to the entity or entities that authorized, implemented, or enabled the execution path consistent with WRS constraints.

In legal terms, the execution boundary functions as a **causal firewall**, ensuring that potential harm remains contained within the speculative domain until and unless it satisfies the WRS permissibility criteria. By structurally separating speculation from action, WRS prevents hypothetical or exploratory reasoning from being misinterpreted as actionable causation.

### 4.5.3 Human Presence Does Not Imply Execution Authority

WRS explicitly rejects the assumption that human involvement automatically confers execution authority. Human-in-the-loop (HITL) mechanisms may provide oversight, validation, or ethical review, but they do not override execution vetoes.

Human presence is a **symbolic safeguard** unless veto power is structurally enforced at the execution boundary. A human who approves an action that violates WRS constraints does not legitimize execution; the system must remain blocked regardless of endorsement.

This design protects human operators and institutions from being placed in positions of illusory control, where nominal approval masks the absence of real execution authority and concentrates legal liability without actual power.

### 4.5.4 Veto Authority Without Decision Ownership

WRS introduces a strict separation between **veto authority** and **decision ownership**. Veto authority governs whether execution may occur; it does not imply authorship, intent, or responsibility for the decision itself.

An execution veto does not assign blame, reverse judgment, or invalidate upstream reasoning. It merely prevents physical commitment. Conversely, permitting execution under WRS does not imply endorsement of the decision's moral, strategic, or social correctness; it indicates only that execution satisfies all non-negotiable constraints.

This neutrality enables consistent governance and avoids retroactive moralization of technical safeguards.

### 4.5.5 Auditability, Evidence Preservation, and Responsibility Anchoring

Because WRS resolves execution permissibility as a Boolean state, accountability can be anchored through auditable records of execution outcomes. Each execution request yields a definitive result—**permitted** or **blocked**—accompanied by a traceable record of which constraints were satisfied or violated.

The Boolean nature of execution ensures **evidence preservation**. Accountability is grounded in verifiable binary logs rather than in subjective interpretations of risk scores, confidence levels, or probabilistic model behavior. This significantly reduces the legal ambiguity often associated with reconstructing intent or causation in AI-driven systems.

Such discrete audit trails support regulatory review, compliance audits, and legal proceedings by providing clear, objective evidence of when and why execution was allowed or denied.

### 4.5.6 Accountability Without Optimization or Delegation

WRS does not delegate accountability to optimization mechanisms, risk models, or adaptive policies. Responsibility cannot be diffused through probabilistic thresholds, automated trade-offs, or emergent system behavior.

Instead, accountability is preserved by enforcing a strict execution boundary: if execution is blocked, no harm occurs; if execution is permitted, responsibility follows the permitted path. This binary structure prevents responsibility dilution and preserves clear lines of causation.

### 4.5.7 Layered Accountability in LERA-Compatible Architectures

Within architectures that incorporate judgment-based governance frameworks such as LERA, WRS enables a model of **layered accountability**.

- **Judgment frameworks** govern deliberative legitimacy: who may decide, under what conditions, and with what authority.

- **WRS** governs execution permissibility: whether physical action may occur at all.

Responsibility is therefore distributed but not diluted. Judgment governs intent and authorization; WRS governs physical consequence. Neither layer substitutes for the other, and neither can override the other's scope.

### 4.5.8 Summary

WRS redefines accountability by anchoring it at the execution boundary. By functioning as a causal firewall, enforcing auditable Boolean outcomes, and separating veto authority from decision ownership, WRS prevents responsibility diffusion while preserving space for intelligent exploration.

This execution-boundary-based accountability model aligns legal responsibility with physical consequence, offering regulators, institutions, and system designers a clear and enforceable framework for governance in the era of autonomous systems.

## 4.6 Applicability Beyond AGI

Although this paper is situated within contemporary discussions of artificial general intelligence (AGI), the applicability of the World Reliability Ruleset (WRS) is neither limited to AGI nor dependent on any specific level of machine intelligence. WRS governs execution rather than cognition, and execution remains invariant across technological paradigms.

### 4.6.1 Intelligence-Agnostic Governance

WRS is explicitly **agnostic to intelligence scale**. Its execution constraints do not reference reasoning capability, model architecture, learning mechanism, or alignment strategy. Whether an action is proposed by a narrow system, an AGI, a human operator, or a future superintelligent system is irrelevant to execution permissibility.

From the perspective of WRS, intelligence affects *what may be proposed*, but never *what may be executed*. This separation ensures that advances in reasoning, prediction, or creativity do not erode the integrity of execution boundaries.

As systems progress from narrow AI to AGI and beyond, the execution boundary remains structurally identical.

### 4.6.2 Compatibility with Artificial Superintelligence (ASI)

The transition from AGI to artificial superintelligence (ASI) amplifies, rather than diminishes, the necessity of WRS. Increased intelligence expands the system's capacity to generate novel strategies, identify unconventional pathways, and justify actions with internally coherent reasoning. None of these properties reduce physical risk.

WRS remains compatible with ASI precisely because it does not depend on trust, alignment confidence, or predictive sufficiency. A system's ability to explain or optimize does not grant it execution privilege.

In this sense, WRS anticipates ASI rather than reacts to it. The ruleset assumes that reasoning capacity may exceed human comprehension and therefore anchors governance at the only invariant interface: execution.

### 4.6.3 Applicability to Non-AI Automated Systems

WRS is equally applicable to automated systems that do not qualify as AI in any contemporary sense. Industrial control systems, autonomous machinery, energy infrastructure, and safety-critical automation already operate at the execution boundary where irreversible consequences occur.

In such systems, failures are rarely cognitive. They are executional. WRS provides a governance layer that is independent of intelligence and therefore suitable for legacy systems, hybrid automation, and emerging cyber-physical infrastructures.

This universality distinguishes WRS from AI-specific governance frameworks and situates it within a broader class of execution governance protocols.

### 4.6.4 Human-Initiated Execution

WRS also applies to execution initiated directly by humans when such execution is mediated through automated or semi-automated systems. Human intent does not bypass execution constraints.

When humans operate within WRS-governed environments, they remain subject to the same execution vetoes as machine-generated actions. This symmetry prevents the creation of privileged execution paths and ensures consistent governance regardless of origin.

Execution authority is therefore detached from authorship, intention, or moral agency and grounded solely in permissibility under WRS constraints.

### 4.6.5 Boundaries of Applicability

WRS does not govern reasoning, learning, communication, simulation, or symbolic manipulation in isolation. Domains such as data processing, internal deliberation, or speculative planning fall outside its scope unless and until they result in execution requests that cross into physical or systemic commitment.

This limitation is intentional. By refusing to regulate cognition, WRS preserves intellectual freedom while enforcing absolute discipline at the point of real-world effect.

### 4.6.6 Execution Governance Is Not a Speed Constraint

A recurring misconception in AI governance is the assumption that stronger execution constraints necessarily imply reduced intelligence velocity or inhibited innovation. WRS explicitly rejects this framing.

WRS does not regulate cognition, learning, optimization, or epistemic exploration. Instead, it provides **functional decoupling between epistemic exploration and kinetic manifestation**, allowing systems to reason, simulate, and explore expansive solution spaces without automatically translating speculative intelligence into physical commitment.

In this sense, WRS treats execution as the **physical constant of governance**, shielding reality from the stochastic volatility introduced by intelligence scaling. As reasoning capacity accelerates—whether in AGI, ASI, or human–machine hybrid systems—the execution boundary remains invariant.

WRS therefore acts not as a brake, but as a **containment vessel**, enabling high-velocity intelligence to operate without the risk of physical spillover. Systems governed by WRS may pursue aggressive optimization, generate extreme hypotheses, or identify unconventional strategies, provided that execution remains subject to non-negotiable constraints.

The execution boundary functions as a **one-way gate**: intelligence may simulate infinite futures, but WRS permits only those trajectories that satisfy execution permissibility to manifest in the physical world. This asymmetry preserves creative freedom while enforcing irreversible safety.

Systems that operate without a formal execution boundary are not faster; they are merely unconstrained. In such systems, execution becomes an implicit continuation of reasoning, allowing probabilistic confidence or optimization success to substitute for physical permissibility. This collapse of boundary, rather than intelligence velocity itself, is the primary source of irreversible harm in automated environments.

WRS does not slow systems down. It prevents them from crossing the only boundary that cannot be reversed.

### 4.6.7 Summary

WRS is not an AGI-specific safeguard but a general execution governance protocol. Its applicability spans narrow automation, AGI, ASI, and human-mediated systems alike. By remaining intelligence-agnostic and execution-focused, WRS provides a stable governance foundation across technological transitions that may fundamentally alter how decisions are generated but not how consequences unfold.

In an era of accelerating intelligence, WRS asserts a simple invariant: execution remains the first irreversible step into risk, and it must therefore remain governed by rules that intelligence cannot override.

## 5. The Boundary of Scope: Non-Goals and Constraints

To preserve the functional integrity and long-term applicability of the World Reliability Ruleset (WRS), it is essential to define the domains it does **not** seek to govern. These non-goals are not limitations of the framework, but structural boundaries that prevent category errors, scope inflation, and misinterpretation.

WRS is designed as an execution governance protocol. Its authority begins and ends at the execution boundary. Any extension beyond this boundary would undermine the very invariance that gives WRS its reliability.

## 5.1 Cognitive Non-Interference

WRS does not regulate cognition.

It does not constrain learning speed, reasoning depth, model capacity, creativity, exploration, or optimization. WRS places no limits on how intelligence is developed, trained, scaled, or deployed at the cognitive level.

This non-interference is deliberate. WRS provides **functional decoupling between epistemic exploration and kinetic manifestation**, allowing systems to reason freely without automatically translating speculative intelligence into physical commitment.

By isolating execution from cognition, WRS expands the safe domain for experimentation. Systems may generate extreme hypotheses, explore unconventional solution spaces, or simulate high-risk scenarios without incurring real-world consequences, provided execution remains subject to WRS constraints.

WRS therefore does not slow intelligence. It shields reality.

## 5.2 Architecture Agnosticism

WRS is agnostic to system architecture.

It does not depend on neural network structures, symbolic reasoning, probabilistic models, agent-based systems, or any specific computational paradigm. Whether intelligence is implemented through transformers, neuromorphic hardware, hybrid symbolic systems, or future quantum architectures is irrelevant to WRS applicability.

This agnosticism ensures that WRS remains valid across technological transitions. As computational substrates evolve, the execution boundary remains invariant.

WRS governs the moment of physical commitment, not the mechanism of thought. As long as execution can produce irreversible physical or systemic consequences, WRS constraints apply.

## 5.3 Separation from Ethics and Value Alignment

WRS is not an ethical framework.

It does not determine what actions are morally right, socially desirable, or politically acceptable. Nor does it encode values, preferences, or normative judgments about outcomes.

WRS governs **physical permissibility**, not moral correctness.

Ethical reasoning, value alignment, and judgment of intent belong to deliberative governance layers, such as LERA. These layers evaluate whether an action *should* be considered. WRS evaluates only whether an action *may* be executed without violating non-negotiable physical constraints.

This separation prevents moral disagreement from weakening execution safety. Regardless of ethical interpretation or value conflict, execution remains blocked if WRS constraints are violated.

## 5.4  No Optimization, No Trade-offs

WRS does not optimize.

It does not assign risk scores, calculate expected utilities, balance competing objectives, or negotiate trade-offs between safety and performance. WRS is explicitly **non-compensatory**: no degree of benefit, confidence, or authorization can offset a violation of execution constraints.

Optimization frameworks operate within solution spaces. WRS defines the boundary of admissibility for those spaces.

By refusing optimization, WRS avoids the gradual erosion of safety that occurs when constraints are treated as tunable parameters rather than absolute conditions.WRS does not guarantee liveness or system availability. A permanently blocked system reflects a governance choice prioritizing irreversibility avoidance over operational continuity.

## 5.5  No Substitution for Responsibility or Judgment

WRS does not replace human responsibility, institutional governance, or legal authority.

It does not decide who is accountable, who is authorized, or who bears moral or legal responsibility. Instead, WRS preserves these structures by preventing execution from occurring in conditions where responsibility would be ambiguous, diffused, or impossible to enforce.

In this sense, **WRS functions as the grounding wire for accountability**. It ensures that the abstract power of intelligence does not discharge uncontrollably into the physical world. By enforcing a stable execution boundary, WRS provides a reliable point at which responsibility can be anchored, audited, and enforced.

WRS neither assumes responsibility nor absolves actors of it. It ensures that when execution does occur, responsibility is transferred across a clearly defined and verifiable boundary.

## 5.6  Summary

The strength of WRS lies not in the breadth of what it governs, but in the precision of what it refuses to govern.

By enforcing cognitive non-interference, architecture agnosticism, ethical separation, non-optimization, and clear responsibility boundaries, WRS maintains its role as an invariant execution boundary rather than an adaptive policy instrument.

**WRS governs not by policy, but by possibility.**
It does not prescribe behavior or outcomes; it defines what is physically admissible. In doing so, WRS proposes a categorical boundary that intelligence—regardless of speed, scale, or sophistication—cannot override.

These non-goals are essential. They protect WRS from misapplication, preserve its neutrality, and ensure that it remains a stable governance primitive across domains, technologies, and eras.

# 6. Conclusion

As artificial intelligence systems advance toward increasingly autonomous and high-velocity operation, the primary governance challenge is no longer how decisions are made, but how-and whether-they are allowed to manifest in the physical world. Contemporary governance approaches have focused extensively on alignment, ethics, and risk assessment, yet they remain structurally incomplete in the absence of a binding execution boundary.

This paper introduced the World Reliability Ruleset (WRS) to address this gap. WRS reframes governance at the point where intelligence becomes consequence. It does not evaluate intent, optimize outcomes, or adjudicate values. Instead, it defines a categorical execution boundary governed by non-negotiable constraints, enforced through deterministic veto logic.

By formalizing execution as a Boolean, non-compensatory state, WRS restores a principle that has been eroded in complex automated systems: **that not all technically possible actions should be physically permissible**. This principle holds regardless of intelligence scale, system architecture, or alignment strategy. As intelligence accelerates—from narrow automation to AGI and beyond—the execution boundary must remain invariant.

WRS deliberately separates cognition from execution, ethics from permissibility, and judgment from enforcement. In doing so, it enables a governance structure in which exploration may scale without destabilizing reality, responsibility may be assigned without ambiguity, and safety may be preserved without constraining innovation. Intelligence is allowed to expand; execution is required to obey.

The implications of this separation extend beyond AI. Any system—human or machine—that can trigger irreversible physical or systemic consequences requires an execution boundary that intelligence cannot override. WRS provides such a boundary, not as policy guidance, but as a technical and normative constant.

This work does not propose a new philosophy of intelligence, nor a comprehensive theory of governance. It proposes a missing primitive. By defining what execution may not do, WRS makes room for intelligence to do more—without transferring speculative power directly into irreversible harm.

In an era where intelligence is increasingly unconstrained, governance must become precise. WRS asserts a simple invariant: **execution is the first irreversible step into risk, and it must therefore remain governed by rules that intelligence cannot negotiate away**.

Future work will explore implementation hardening strategies, emergency governance interfaces, and the formal mapping of WRS constraints onto legal liability and regulatory audit frameworks, while preserving the core execution semantics proposed in this framework.

# References

[1] J. L. Liu, "LERA: Reinstating Judgment as a Structural Precondition for Execution in Automated Systems," arXiv preprint arXiv:2601.08880 [cs.CY], 2026.

[2] J. L. Liu, "The World Reliability Ruleset (WRS): A Technical Specification Supporting the Structural Execution Boundary Framework," Technical Report, 2026.

[3] L. Floridi et al., "AI4People—An Ethical Framework for a Good AI Society," Minds and Machines, vol. 28, no. 4, pp. 689–707, 2018.

[4] D. Amodei et al., "Concrete Problems in AI Safety," arXiv preprint arXiv:1606.06565, 2016.

[5] NIST, AI Risk Management Framework (AI RMF 1.0), National Institute of Standards and Technology, 2023.

[6] IEEE Standards Association, IEEE P7000 Series: Model Process for Addressing Ethical Concerns During System Design, 2020.

[7] Future of Life Institute, Asilomar AI Principles, 2017.

[8] Y. Bai et al., "Constitutional AI: Harmlessness from AI Feedback," arXiv preprint arXiv:2212.08073, 2022.

[9] C. Perrow, Normal Accidents: Living with High-Risk Technologies. Princeton, NJ, USA: Princeton University Press, 1984.

[10] U. Beck, Risk Society: Towards a New Modernity. London, UK: Sage Publications, 1992.

[11] N. Leveson, Engineering a Safer World: Systems Thinking Applied to Safety. Cambridge, MA, USA: MIT Press, 2011.

# Appendix A: Interoperability with Existing Governance Frameworks

This appendix clarifies how the World Reliability Ruleset (WRS) interoperates with existing AI governance, risk management, and alignment frameworks. Its purpose is not comparative evaluation, but **technical realization**: to specify how WRS functions as a **deterministic enforcement layer** that converts governance intent into executable constraints.

WRS is designed to be **orthogonal** to decision-centric, ethics-centric, and risk-centric frameworks. It does not replace them. It provides the execution-layer mechanism required for their reliable operation in systems with irreversible consequences.

## A.1 WRS as a Deterministic Enforcement Layer

Across contemporary governance regimes, a consistent structural limitation can be observed: governance requirements are defined normatively, but enforced procedurally or retrospectively.

WRS addresses this limitation by acting as a **technical realization of governance constraints**. It transforms static legal and policy requirements into **dynamic, deterministic execution control**.

In this role, WRS functions as a **deterministic enforcement layer**—a non-negotiable execution boundary that operates independently of interpretive judgment, organizational discretion, or probabilistic risk scoring.

WRS does not redefine governance goals. It ensures that once defined, those goals cannot be violated at the moment of execution.

## A.2 Interoperability Overview

### Table 2. Interoperability Between WRS and Existing AI Governance Frameworks

| Governance Framework | Primary Governance Focus | Structural Limitation | WRS Interoperability Role |
|---|---|---|---|
| NIST AI Risk Management Framework (AI RMF) | Risk identification and mitigation | Procedural rather than enforceable controls | WRS provides deterministic execution enforcement beneath RMF processes |
| EU AI Act | Legal classification and compliance obligations | Limited pre-execution technical blocking | WRS enables real-time execution vetoes for high-risk systems |
| Constitutional AI (Anthropic) | Normative alignment of model behavior | No physical execution control | WRS introduces a final Boolean execution gate |
| IEEE P7000 Series [6] | Ethical system design standards | Ethics without hard execution primitives | WRS translates ethical intent into enforceable execution constraints |

*Table 2 illustrates how WRS interoperates with existing AI governance frameworks by providing a deterministic execution enforcement layer that complements decision-, risk-,*

*ethics-, and compliance-centric approaches.This interoperability is **non-hierarchical**. WRS does not subsume these frameworks, nor is it subsumed by them. It provides the missing execution boundary that allows their governance objectives to be realized in practice.*

## A.3 WRS and the NIST AI Risk Management Framework

The NIST AI RMF [5] emphasizes organizational processes for identifying, measuring, and mitigating AI-related risks. While effective for governance planning, RMF relies on human interpretation and procedural compliance.

WRS complements RMF by providing **deterministic enforcement at the execution layer**. Risk assessments may inform decisions, but execution under WRS remains subject to absolute permissibility constraints that cannot be overridden by confidence, urgency, or institutional pressure.

In effect, WRS transforms RMF from a risk documentation framework into an **operationally enforceable system**.

## A.4 WRS and the EU AI Act

The EU AI Act establishes legal obligations for high-risk AI systems, focusing on transparency, oversight, and accountability. However, compliance mechanisms are largely procedural and post hoc.

WRS provides a **technical realization of these legal obligations** by enabling pre-execution blocking of impermissible actions. Rather than relying on after-the-fact enforcement, WRS ensures that prohibited execution paths are structurally inaccessible.

In this role, WRS converts legal compliance from static documentation into **dynamic technical enforcement**.

## A.5 WRS and Constitutional AI

Constitutional AI approaches guide model behavior through internal normative principles. These frameworks operate entirely within the cognitive domain.

WRS does not interfere with constitutional reasoning. Instead, it enforces a **final, non-negotiable Boolean execution gate** after all internal reasoning and alignment processes have concluded.

This separation preserves normative flexibility while ensuring that irreversible physical execution remains governed by deterministic constraints.

## A.6 Orthogonality and Non-Substitutability

WRS is not a substitute for governance frameworks, ethical principles, or legal standards. It is the execution-layer condition that enables them to function as intended.

Governance frameworks define *what should be done*.

WRS defines *what can be executed*.

This orthogonality ensures that WRS can be integrated into diverse regulatory, organizational, and philosophical contexts without requiring consensus on values or objectives.

## A.7 Summary

Existing governance frameworks articulate intent.WRS provides enforcement.

**Governance without execution vetoes is merely documentation.WRS supplies the deterministic execution control—the "teeth"—that existing frameworks currently lack.**

By transforming governance requirements into enforceable execution boundaries, WRS enables legal, ethical, and risk-based frameworks to achieve operational closure in systems where execution entails irreversible physical or systemic consequences.

WRS is therefore not an alternative governance model, but the **technical foundation that makes governance real**.

## Appendix B: Rule Structure and Trigger Semantics of WRS

This appendix formalizes the rule structure underlying the World Reliability Ruleset (WRS) and clarifies how execution constraints are triggered, evaluated, and enforced. The purpose of this section is not to enumerate domain-specific safety parameters, but to define the logical form and enforcement semantics that qualify WRS as a ruleset [2] rather than a policy, guideline, or heuristic framework.Any system capable of rewriting or bypassing its own execution boundary is, by definition, operating outside the scope of WRS governance.

### B.1 Rule Form and Trigger Phase

In WRS, rules are expressed as *execution constraints* rather than recommendations, objectives, or probabilistic thresholds. Each constraint specifies a non-negotiable condition that must be satisfied for an execution event to be permitted.

Constraints are evaluated at the **pre-commitment phase** of execution, prior to any irreversible physical, kinetic, or systemic state transition. This ensures that the transformation from abstract logic to physical manifestation is gated by the ruleset itself, rather than assessed retrospectively. The evaluation occurs at the execution boundary, independent of decision quality, intent, optimization performance, or predicted utility.

By locating enforcement at this pre-commitment boundary, WRS ensures that execution control remains temporally and causally prior to physical commitment, even in systems operating at high speed or fine temporal resolution.

WRS does not assume software-level evaluation latency. The execution boundary may be realized through hardware interlocks, firmware-level gates, or physically enforced control circuits whose response time precedes kinetic commitment.

### B.2 Boolean Trigger Semantics

All WRS constraints operate under a Boolean trigger semantics. At the moment of execution, each constraint is evaluated as either satisfied (True) or violated (False). No intermediate confidence scores, weights, or probabilistic gradients are permitted within the core enforcement logic.

Formally, let

$$C = \{c_1, c_2, \ldots c_n\}$$

denote the set of applicable execution constraints. Execution is permitted if and only if:

$$\bigwedge_{i=1}^{n} eval(c_i) = True$$

Any single violated constraint is sufficient to block execution.

## B.3 Veto Logic and Non-Compensatory Enforcement

WRS adopts a *single-veto sufficiency* principle. Constraint violations are non-compensatory: no combination of satisfied constraints may offset a violated one. This distinguishes WRS from optimization-based or risk-balancing governance approaches in which trade-offs are allowed.

The veto mechanism is absolute with respect to the execution event. Once triggered, execution is categorically prohibited regardless of system confidence, administrative authorization, or anticipated benefits.

## B.4 Core and Domain-Specific Rule Sets

WRS consists of a core ruleset (WRS-C) and optional domain-specific constraint subsets (WRS-$D_x$). WRS-C defines invariant execution principles applicable to all systems with irreversible consequences. WRS-D subsets specialize these principles for particular physical domains without altering the underlying veto semantics.

Domain-specific constraints may introduce additional blocking conditions but may not weaken, override, or bypass any core constraint. Constraint inheritance is strictly one-directional: all WRS-D executions remain subject to WRS-C veto authority.

A domain-specific subset is derived whenever execution within a system introduces a distinct class of irreversible physical or systemic risk that cannot be fully governed by WRS-C alone.

## B.5 Auditability, Attribution, and Extensibility

The Boolean trigger semantics of WRS ensure that every veto event is **uniquely attributable** to one or more violated constraints. This enables deterministic post-hoc auditing, eliminating ambiguity in responsibility attribution and preventing retrospective reinterpretation of probabilistic risk scores.

The set of domain-specific constraint subsets defined in this paper is non-exhaustive. The absence of a domain-specific subset does not imply exemption from WRS applicability. Any system whose execution entails irreversible physical, kinetic, or systemic consequences falls within the scope of WRS, regardless of intelligence level, implementation architecture, or deployment context.