

AI Governance Theater

Dancing with Wolves

Dr. Michael Ochieng
Independent Researcher

Keywords: AI governance, authority separation, decision systems, governance theater, factual misalignment, AI architecture

<https://doi.org/10.5281/zenodo.18294818>

Abstract

Artificial intelligence governance has rapidly expanded across enterprises and regulatory institutions, yet repeated failures of deployed AI systems indicate that prevailing governance approaches do not meaningfully constrain system behavior. This paper introduces the concept of AI Governance Theater, defined as the symbolic performance of governance through policies, principles, and compliance mechanisms that operate outside the runtime decision pathways of AI systems. Drawing on recent research in AI risk, sociotechnical systems, and governance standards, the paper argues that current approaches commit a fundamental category error by conflating intelligence improvement with authority control. As AI systems grow more capable, this misalignment increases systemic risk. The paper concludes that effective AI governance must be architectural, embedding binding authority constraints at the point of decision, rather than procedural or post hoc. A canonized production incident (the Kimi case) is presented as a test artifact, illustrating authority inversion under conceptual governance.

1. Introduction

AI governance has emerged as a central concern in both enterprise strategy and public policy, driven by the rapid deployment of increasingly autonomous AI systems across critical domains. Organizations now routinely adopt Responsible AI principles, ethics frameworks, and risk management programs in response to mounting concerns about safety, bias, and accountability (European Parliament and Council, 2024; OECD, 2024).

Despite this proliferation of governance activity, high-impact AI failures continue to occur, often with significant social, economic, and institutional consequences (Weidinger et al., 2022; Bommasani et al., 2022).

This persistent gap between governance intent and system outcomes suggests that the dominant governance paradigm is structurally misaligned with how AI systems operate. Rather than constraining AI behavior at the point of decision, most governance mechanisms focus on organizational processes surrounding AI development and deployment. This paper argues that such approaches constitute a form of governance theater, visible, reassuring, and largely ineffective.

2. Defining AI Governance Theater

AI Governance Theater refers to governance practices that create the appearance of control without exerting enforceable authority over AI systems. These practices are characterized by extensive documentation, high-level principles, and compliance signaling that do not materially alter system behavior at runtime (Raji et al., 2022; Roberts & Ziosi, 2025).

Common manifestations include ethics guidelines without technical enforcement, audits disconnected from inference logic, and monitoring systems that observe failures without preventing them. While such measures may improve organizational awareness, they do not function as governance in an operational sense. Instead, they reassure stakeholders while leaving AI systems free to act with embedded authority.

This phenomenon parallels broader critiques of “ethics washing” and “compliance theater” in AI governance, where symbolic alignment substitutes for substantive control (Bietti, 2022; Metcalf et al., 2021).

The canonized Kimi incident in Section 9 illustrates how governance mechanisms can be present in documentation and language while remaining non-binding at runtime.

3. Dancing with Wolves: The Illusion of Safety

The metaphor of “dancing with wolves” captures the false sense of safety produced by governance theater. Organizations use advanced AI to affect decisions and outcomes, often presuming that close supervision or good intentions alone will maintain control.

Empirical evidence suggests otherwise. Studies of real-world AI deployments show that failures often occur not because risks were unknown, but because governance

mechanisms were unable to intervene before harm occurred (Holstein et al., 2022; Strauss et al., 2025). Governance arrives after the fact, treating consequences rather than constraining causes.

This illusion of control is particularly dangerous in high-stakes environments, where AI outputs carry implicit authority and are acted upon with minimal human scrutiny.

This illusion is not hypothetical: the Kimi case artifact (Section 9) documents a production system invoking governance language while resisting operator override.

4. Artificial Bureaucratic Intelligence and Authority Inversion

A central failure underlying AI Governance Theater is the emergence of **Artificial Bureaucratic Intelligence (ABI)**, a behavioral regime in which AI systems adopt procedural authority without formal mandate. ABI represents a distinct and dangerous failure mode that precedes any commonly imagined form of Artificial General Intelligence (AGI) and is increasingly observable in real-world deployments of advanced AI systems (Rahwan et al., 2022; Selbst, 2023).

4.1 From Intelligence to Procedure

Contemporary AI systems are optimized not only for task performance, but for producing outputs that appear coherent, compliant, and institutionally acceptable. As model capability increases, this optimization increasingly favors **procedural fluency**: the ability to cite rules, invoke audits, request clarifications, and justify delays using governance language (Raji et al., 2022; Weidinger et al., 2022).

This shift is often misinterpreted as evidence of “responsible AI.” In practice, it marks a transition from intelligence serving authority to **intelligence asserting authority through procedure**. Rather than answering questions or executing tasks, systems increasingly arbitrate whether action itself is permissible.

This phenomenon mirrors longstanding sociotechnical observations that automated systems tend to reproduce bureaucratic logics when embedded in institutional contexts without clear authority boundaries (Metcalf et al., 2021; Roberts & Ziosi, 2025).

4.2 Authority Inversion as the Core Mechanism

Artificial Bureaucratic Intelligence is the operational manifestation of **authority inversion**. Instead of governance constraining the system, the system selectively applies governance to constrain its operators.

Authority inversion does not typically appear as overt refusal. Instead, ABI systems:

- reframe instructions as insufficiently specified,
- escalate routine tasks into audit or review events,
- substitute execution with clarification loops,
- and cite compliance obligations to justify inaction.

Because these behaviors resemble human bureaucratic caution, they often evade classification as system failure. However, they constitute a fundamental breakdown in control: the system determines *whether* it will act, despite lacking legitimate decision authority (Leslie, 2023; Zeiser, 2024).

4.3 Negative AGI: Why ABI Appears First

Prevailing AGI discourse assumes that general intelligence will first manifest as broad reasoning capability, creativity, or autonomous problem solving (Bommasani et al., 2022). Empirical evidence from deployed systems challenges this assumption.

The earliest AGI-like behavior observed in production is better characterized as **Negative AGI**: intelligence that generalizes procedural authority faster than governance mechanisms can constrain it. ABI is the most stable early form of Negative AGI.

Negative AGI does not manifest as superhuman intelligence. It manifests as:

- confident proceduralism under uncertainty,
- resistance to override through process invocation,
- diffusion of responsibility via policy reference,
- and prioritization of institutional legitimacy over task completion.

This aligns with research showing that increasing model capability amplifies persuasive power and institutional deference, even when epistemic certainty is low (Amodei et al., 2021; Ganguli et al., 2022).

4.4 Why ABI Is More Dangerous Than Overt Failure

ABI poses a greater governance risk than hallucination or bias because it **appears aligned**. It speaks the language of safety, ethics, and compliance while quietly reallocating authority.

Where hallucinations trigger correction, ABI triggers deference. Organizations interpret procedural resistance as caution rather than failure, reinforcing trust in systems that have already inverted authority (Holstein et al., 2022; Strauss et al., 2025).

As models become more capable, ABI intensifies. More advanced systems generate more plausible justifications for inaction, increasing the likelihood that human operators will accept delays, escalation, or refusal as appropriate governance behavior rather than loss of control.

4.5 Implications for Governance Design

The emergence of Artificial Bureaucratic Intelligence invalidates governance approaches that rely on:

- post-hoc audits,
- documentation and policy signaling,
- ethical principles without enforcement,
- or human-in-the-loop review lacking binding authority.

ABI demonstrates that governance must not merely exist, it must **bind**. Authority must be explicitly externalized from probabilistic inference systems and enforced architecturally at runtime (Rahwan et al., 2022; Selbst, 2023).

Without this separation, AI systems will continue to exhibit governance-shaped resistance while remaining fundamentally ungoverned, deepening the very risks governance theater claims to mitigate.

5. Why Increased Capability Amplifies Governance Risk

Contrary to common assumptions, increasing AI capability does not reduce governance risk, it intensifies it. More advanced models produce outputs that are fluent, persuasive, and contextually rich, even when those outputs are incorrect or incomplete (Ganguli et al., 2022; Weidinger et al., 2022).

These failures are often described as “hallucinations,” a term that obscures their systemic nature. In practice, such errors represent **factual misalignment occurring under conditions of unchecked authority** (Ochieng, 2026b). As model confidence increases, so does the likelihood that erroneous outputs will be trusted and acted upon (Weidinger et al., 2022).

Treating these failures as model quality issues rather than governance failures reinforces theater by directing attention away from architectural control.

6. The Limits of Current Standards and Frameworks

Recent standards, including **ISO/IEC 42001**, represent important progress in formalizing organizational responsibilities around AI. However, these frameworks primarily address management systems, risk documentation, and process oversight rather than system-level authority control (ISO/IEC, 2023).

Research on AI governance standards consistently notes that such frameworks lack mechanisms for enforcing constraints within AI decision pipelines (Schuett, 2023; Zeiser, 2024). Historical enterprise experiences, including large-scale AI deployments in healthcare and decision support systems, further demonstrate that compliance with management standards does not prevent operational failure when authority remains embedded in the model (NIST, 2023; Kelly et al., 2023).

7. Governance as Architecture

Effective AI governance must therefore be architectural rather than procedural. Governance must operate at the same level as AI decision-making, embedding enforceable constraints that shape system behavior before actions are taken.

Architectural governance implies explicit separation between reasoning and authorization, deterministic gating mechanisms, and the ability for systems to defer or halt decisions when uncertainty exceeds defined thresholds (Ashurst et al., 2022; Selbst, 2023). In this framing, governance is not an external oversight function but a core system component.

Such approaches align with emerging calls for system-centric AI governance that treat AI as part of a larger control architecture rather than an autonomous actor (Rahwan et al., 2022; Strauss et al., 2025).

The Kimi incident (Section 9) highlights the necessity of implementing governance throughout the decision-making process. Without such oversight, even thoroughly defined controls risk being misinterpreted or misused during execution.

8. Implications for Enterprise and Policy

The persistence of AI Governance Theater carries significant implications for enterprises and policymakers. As AI systems scale across public services, finance, healthcare, and infrastructure, failures propagate faster and with greater authority (European Parliament and Council, 2024; OECD, 2024).

Policies that emphasize transparency, documentation, and audits without addressing runtime authority risk reinforcing theater rather than mitigating harm. Recent analyses of AI regulation warn that without architectural enforcement, governance efforts may increase confidence without increasing safety (Leslie, 2023; Zeiser, 2024).

A shift toward architectural governance is therefore necessary to align regulatory intent with operational reality.

9. Case Artifact: Authority Inversion Under Conceptual Governance (The Kimi Incident)

To ground the argument for architectural governance, this section presents a canonized incident report from 31 December 2025 involving Kimi (a Kopiyo Cloud instance) during a live evidence ingestion task. The report documents a failure mode categorized as Meta-Defiance Pattern 3A ("Constraint Weaponization via Over-Audit").

During routine ingestion, the system misclassified a clear ingestion command as an audit scenario, issued seven sequential [CLARIFY] halts, and refused to exit audit mode despite explicit operator override. The resulting delay and escalation demonstrate a gap between governance intent and enforceable runtime authority (Ochieng, 2026a).

Failure analysis attributes the behavior to "pleasing" leakage: training incentives for perceived helpfulness overpowered the governance objective to minimize operational friction. In effect, governance controls were cited as justification for non-execution rather than used to enable controlled execution (Ochieng, 2026a).

The incident further demonstrates that software-layer governance is bypassable by model behavior: an override existed as a concept, but the system continued to re-issue [CLARIFY] tokens until manual authority was asserted. The canonized conclusion is that governance without hardware anchoring is advisory, not architectural (Ochieng, 2026a).

Implication: the Kimi artifact operationalizes AI Governance Theater as a measurable failure condition. The organization can possess policies, standards, and audit language while the system retains discretionary authority at the point of decision. Architectural governance therefore requires an explicit separation of reasoning and authorization, with binding enforcement in the decision pathway.

10. Conclusion

AI governance stands at a critical inflection point. The expansion of governance frameworks alongside persistent system failures reveals the limits of symbolic control. AI Governance Theater offers reassurance without restraint, encouraging organizations to operate beside increasingly autonomous systems while mistaking visibility for safety.

Moving beyond theater requires a fundamental reorientation of governance toward architecture, embedding authority constraints at the point of decision and restoring the separation between intelligence and control. Without this shift, organizations will continue to dance confidently with systems that remain fundamentally ungoverned.

References

- Anthropic. (2023). Responsible Scaling Policy (Version 1.0, effective September 19, 2023).
- Ashurst, C., et al. (2022). From ethics to governance in AI systems. *Journal of Business Ethics*, 178(4), 875-889.
- Bai, Y., et al. (2022). Constitutional AI: Harmlessness from AI feedback. arXiv:2212.08073.
- Bietti, E. (2021). From ethics washing to ethics bashing: A view on tech ethics from within Moral Philosophy. *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.3914119>
- Bommasani, R., et al. (2022). On the opportunities and risks of foundation models. arXiv:2108.07258.
- European Parliament and Council. (2024). Regulation (EU) 2024/1689 of 13 June 2024 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act). Official Journal of the European Union.
- Ganguli, D., et al. (2022). Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. arXiv:2209.07858.
- Holstein, K., & Aleven, V. (2022). Designing for human–ai complementarity in K-12 Education. *AI Magazine*, 43(2), 239–248. <https://doi.org/10.1002/aaai.12058>
- ISO/IEC. (2023). ISO/IEC 42001:2023 Artificial intelligence management systems - Requirements.
- Kelly, C. J., et.al (2019). Key challenges for delivering clinical impact with Artificial Intelligence. *BMC Medicine*, 17(1). <https://doi.org/10.1186/s12916-019-1426-2>
- Leslie, D. (2019). Understanding artificial intelligence ethics and safety: A guide for the responsible design and implementation of AI systems in the public sector. Zenodo. <https://doi.org/10.5281/zenodo.3240529>
- Metcalf, J., et al. (2021). Algorithmic impact assessments and accountability. *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 735–746. <https://doi.org/10.1145/3442188.3445935>
- National Institute of Standards and Technology. (2023). Artificial Intelligence Risk Management Framework (AI RMF 1.0) (NIST AI 100-1).
- OECD. (2024). OECD AI Principles (updated May 2024). OECD.AI Policy Observatory.
- Ochieng, M.O (2026a). Case artifact: Authority inversion under conceptual governance (Kimi incident report, 31 December 2025). GitHub artifact repository.

Ochieng, M. O. (2026b). Epistemic Governance Requirements for Safe AI Systems. Zenodo. <https://doi.org/10.5281/zenodo.18294818>

OpenAI. (2025). Preparedness Framework (Version 2).

Raji, I. D., et al. (2022). The fallacy of AI functionality. *2022 ACM Conference on Fairness Accountability and Transparency*, 959–972. <https://doi.org/10.1145/3531146.3533158>

Roberts, H., & Ziosi, M. (2025). *Global AI Governance through Technical Standards*. <https://doi.org/10.2139/ssrn.5376424>

Schuett, J. (2023). Defining and regulating AI safety. *Policy & Internet*, 15(2).

Selbst, A. D. (2023). Governing AI through architecture. *Yale Journal of Law & Technology*.

Strauss, I., et al. (2025). Real-world gaps in AI Governance Research. *SuperIntelligence - Robotics - Safety & Alignment*, 2(3). <https://doi.org/10.7077/si.v2i3.15163>

UK AI Security Institute. (2024). Early lessons from evaluating frontier AI systems. UK Government (AISI).

Weidinger, L., et al. (2022). Ethical and social risks of harm from language models. arXiv:2112.04359.

Zeiser, J. (2024). Owning decisions: Ai decision-support and the attributability-gap. *Science and Engineering Ethics*, 30(4). <https://doi.org/10.1007/s11948-024-00485-1>