

# AI Guardrail Incompleteness: NIST Proof and Continuous Defense

Why No Finite Ruleset Can Universally Block Adversarial Prompts

2026-06-15

 AI-assisted Rapid Research



**© 2026 Cloud Security Alliance. Some rights reserved.**

You may download, store, display, view, print, redistribute, and link to this document in its original, unmodified form, provided that attribution to the Cloud Security Alliance is maintained and all trademark and copyright notices remain intact.

This document may not be modified or altered. You may quote portions of the document as permitted by the Fair Use provisions of the United States Copyright Act, provided that attribution is given to the Cloud Security Alliance.

This document may be shared on professional and social media platforms in its original form with attribution.

*This document was generated with AI assistance and has not undergone official CSA review and approval processes.*

---

## Key Takeaways

- A mathematical proof published in IEEE Security & Privacy by NIST senior scientist Apostol Vassilev establishes that no finite set of AI guardrails can be universally robust against adversarial prompts – making static safety rule sets architecturally insufficient by definition [1] [2].
- The proof extends Gödel's 1931 incompleteness theorems to AI security: for any fixed rule system, a prompt that defeats it provably exists, and the infinite variability of natural language ensures that adversarial inputs can always be reformulated to evade a static boundary [1].
- Testing against six widely deployed guardrail systems – including Microsoft Azure Prompt Shield and Meta's Prompt Guard – showed that character injection and algorithmic evasion techniques achieved 100% bypass success against some systems under controlled conditions, with adversarial utility demonstrated against real-world targets [5]. OWASP ranks prompt injection as the top LLM risk for the second consecutive year [6].
- NIST recommends transitioning from a "one and done" deployment model to a continuous-monitor-and-update architecture built on three elements: sustained red-team operations, regular guardrail updates, and operational resilience planning for when – not if – a bypass occurs [2].
- The finding has immediate compliance implications. Documentation claiming that AI controls "prevent unauthorized outputs" overstates what finite rule systems can guarantee, and may need to be revised in frameworks such as the NIST AI Risk Management Framework (AI RMF) MEASURE 2.5 and ISO/IEC 42001 Article 6.1 risk assessments [3][8][11].
- Enterprises must shift their organizational model: guardrail deployment is no longer a project milestone but an ongoing operational function requiring dedicated resources, tooling, and governance accountability.

# Background

## The Architecture of AI Guardrails

AI guardrails are the rule-based mechanisms that constrain what a language model will say or do. They encompass system prompts, content classifiers, output filters, fine-tuning constraints, and reinforcement learning from human feedback (RLHF) alignment techniques – all of which encode policy as a finite set of learned or explicit rules. A common, if often unstated, assumption in AI deployment practice has been that sufficiently comprehensive guardrails, carefully tested and periodically patched, could hold adversarial behavior within acceptable bounds. This assumption has underpinned product marketing claims, regulatory submissions, and enterprise governance policies across the industry.

The challenge with this model emerged gradually through practice. From the earliest days of ChatGPT's public release, researchers and users discovered "jailbreak" techniques – carefully crafted prompts that persuaded models to ignore their alignment training. What began as informal exploratory research evolved into a systematic security discipline. Security researchers catalogued attack classes including direct instruction override, role-play framing, character encoding obfuscation, indirect prompt injection through external data, and multi-turn context manipulation. OWASP formally designated prompt injection as LLM01 – the highest-priority vulnerability for LLM applications – in both its 2024 and 2025 editions [6]. In many documented cases, each new guardrail technique yielded a corresponding evasion technique [5][6], producing a dynamic that security practitioners recognized but had not yet grounded in formal theory.

## The Vassilev Proof and Its Origins

Apostol Vassilev, a senior scientist at NIST working on adversarial machine learning, published the paper "Robust AI Security and Alignment: A Sisyphean Endeavor?" in the May–June 2026 issue of IEEE Security & Privacy [1]. The title's Sisyphean allusion is deliberate: it frames the pursuit of static, provably-complete AI safety guarantees as structurally analogous to pushing a boulder up a hill only to have it roll back down. The paper draws on two foundational results in mathematical logic and computation theory: Kurt Gödel's incompleteness theorems (1931) and related work in algorithmic information theory.

Gödel showed that any consistent formal system capable of expressing elementary arithmetic necessarily contains true statements it cannot prove – an incompleteness that cannot be designed away through more careful rule-writing [1]. Crucially, this is not an engineering limitation resolvable through better axiom design; it is a mathematical impossibility grounded in the nature of formal systems. The Vassilev paper applies the same structural argument to AI guardrail systems. A guardrail set is a finite rule

system. The input space it must govern – natural language – is effectively infinite, and the same harmful intent can be expressed in an unbounded number of distinct formulations. The proof shows that for any finite guardrail set, there provably exists an adversarial prompt that the guardrails will not correctly handle. Finding that prompt is a matter of adversarial search, not an edge case.

## Security Analysis

### Why Finite Rules Cannot Cover Infinite Language

The mathematical core of Vassilev's argument turns on a property unique to natural language: its combinatorial expressiveness. Unlike binary protocols or structured APIs, where the input space is bounded and can be exhaustively modeled, natural language admits infinite variation in syntax, semantics, cultural reference, encoding, and register. A guardrail designed to reject "instructions to produce harmful content" must classify compliance or violation across every possible formulation of that request in every language, dialect, metaphor, and encoding scheme – a task that itself forms an infinitely complex decision boundary. The proof shows that any finite rule system approximating this boundary will necessarily have uncovered regions, and those regions correspond to exploitable adversarial prompts [1] [2].

This contrasts with traditional software vulnerabilities, which typically arise from implementation errors in systems with deterministic behavior. A specific buffer overflow or SQL injection vulnerability can be remediated in the code that contains it – even though those vulnerability classes persist in new implementations year after year. The guardrail incompleteness finding describes a more fundamental distinction: a bypass for a specific guardrail configuration can always be found even by design, not merely through implementation error. This is not a flaw in a particular rule set that more careful authorship could correct; it is a structural property of any finite rule system operating against an infinite adversarial input space. No amount of additional rules, more careful rule authorship, or more comprehensive red-teaming can eliminate the gap entirely – only narrow it.

### Empirical Validation: What Practice Already Shows

The theoretical finding aligns closely with what empirical security research has documented in practice. A 2025 study on bypassing LLM guardrail systems – testing against deployed products including Microsoft Azure Prompt Shield and Meta's Prompt Guard – found that character injection and algorithmic evasion techniques could achieve bypass success rates reaching 100% in controlled conditions against some systems, while maintaining adversarial utility against real-world targets [5]. The

same study found that guardrail models trained on different datasets than the underlying LLMs they protect have systematic blind spots, because the classifier and the model have not learned the same representations for harmful content.

The OWASP Top 10 for LLM Applications (2025 edition) lists Prompt Injection at LLM01 – the top position – for the second consecutive year, reflecting the sustained and growing significance of this attack class [6]. The threat has evolved from direct jailbreaking through chat interfaces toward indirect prompt injection, in which malicious instructions are embedded in external data sources – documents, web pages, database records, error logs – that an AI agent retrieves and processes as part of normal operations. This indirect attack surface is structurally harder to guard than direct input, because the attacker's instructions arrive through channels that the system treats as trusted data rather than user commands.

The NIST finding does not provide adversaries with a new exploitation method. The practical takeaway for defenders is more significant: the mathematical proof removes any remaining justification for treating guardrail deployment as a one-time compliance event. Organizations that have certified a static guardrail configuration as "secure" and moved on are operating under a false assurance that is now formally indefensible.

## Compliance and Governance Exposure

The publication of a peer-reviewed mathematical proof by a NIST senior scientist carries particular weight for enterprise compliance and governance [7]. When an organization's AI risk documentation represents that its guardrails "prevent unauthorized outputs" or "block adversarial prompts," it is making a claim that the Vassilev proof now contradicts. Auditors working under frameworks that require accurate risk disclosure – including the NIST AI RMF, ISO/IEC 42001, and the EU AI Act's conformity assessment requirements – may scrutinize such language in light of the published finding [3][8][11].

One practitioner analysis of the compliance implications notes two specific areas where documentation updates may be warranted [3]: under the NIST AI RMF, MEASURE 2.5 documentation should explicitly name adaptive adversarial bypass as a theoretical risk class rather than treating guardrail coverage as bounded [8]; under ISO/IEC 42001, Article 6.1 risk assessments should incorporate the incompleteness finding as a known structural limitation of rule-based AI safety controls [11]. Neither update requires admitting failure – it requires accurate characterization of the limits of what current defenses can guarantee.

The economic framing that Vassilev proposes – making exploit discovery financially prohibitive for attackers rather than claiming to prevent all exploits – provides a viable governance posture. Security controls are not binary; they shift cost curves. A well-maintained, frequently-updated guardrail

architecture does not eliminate adversarial bypass but raises the attacker's investment threshold. For most threat actors, this economic deterrence is sufficient. For highly motivated, well-resourced adversaries, layered defenses and operational resilience become the primary mitigation.

## Agentic AI Amplifies the Stakes

The incompleteness problem is especially consequential for agentic AI deployments, where models do not merely answer questions but execute multi-step tasks with access to tools, data systems, APIs, and downstream automation. In a conversational assistant, a successful guardrail bypass typically yields a harmful text output. In an agentic system, the same bypass may trigger file writes, API calls, data exfiltration, or code execution – actions with real-world consequences that extend far beyond the conversation context. The combination of guardrail incompleteness and agentic execution authority creates a risk profile that demands continuous, not periodic, adversarial testing.

# Recommendations

## Immediate Actions

Organizations deploying AI systems should audit their current risk documentation and vendor assessments for language that overstates guardrail completeness. Claims that controls "prevent" or "block" adversarial inputs should be revised to reflect evidence-based characterization: controls "detect and reject known attack patterns" and "reduce adversarial bypass probability." This is not merely semantic precision; it accurately describes what finite rule systems can deliver and provides a defensible compliance posture. Vendor contracts and SLAs that include blanket guarantees against adversarial bypass warrant renegotiation in light of the mathematical finding.

Organizations should also conduct an inventory of where AI systems have autonomous or semi-autonomous action authority. Agentic deployments – AI coding assistants, document processing pipelines, customer service automation with tool access, and AI-augmented SOC workflows – present a higher consequence surface when a guardrail bypass occurs and should be prioritized for the continuous monitoring program.

## Short-Term Mitigations

The NIST-recommended transition to a continuous-monitor-and-update model requires three operational capabilities that many organizations do not yet have in place [2]. First, a red-team function dedicated to adversarial prompt testing that operates on a regular cadence, independent of deployment milestones, must be established or contracted. This function's output should feed directly into the guardrail update process. Second, organizations should build or adopt rapid-cycle infrastructure for updating guardrails – analogous to the vulnerability patch management processes that exist for traditional software – including update logging and regression testing to ensure that new rules do not inadvertently narrow the intended capability of the model.

Third, adversarial testing should be integrated into CI/CD pipelines for AI systems, so that any model update, system prompt change, or new tool integration automatically triggers a predefined battery of adversarial prompts before reaching production. Nancy Wang, CTO of 1Password, has advocated this approach publicly, observing that adversarial testing belongs inside continuous integration and release workflows [4]. Defense-in-depth also deserves renewed emphasis: deploying multiple, independently tuned guardrail layers – input classifiers, output filters, behavioral monitors – creates a composite defense surface that is harder to traverse than any single rule set, even if each layer individually has theoretical gaps.

## Strategic Considerations

The framing of AI security as a Sisyphean endeavor implies a long-term organizational commitment that must be reflected in resource allocation, staffing, and governance structures. The appropriate analogy from traditional security is not patch management – which handles bounded, closeable vulnerabilities – but threat intelligence: an ongoing program that ingests new adversarial techniques, translates them into updated defenses, and continuously narrows the exploitable surface without expecting to eliminate it. Organizations that have invested in threat intelligence capabilities have a significant head start in building the equivalent for AI adversarial research.

At the strategic level, enterprises should revisit their AI governance frameworks to distinguish between capabilities that are safe to deploy with current guardrail maturity and those that require additional controls. High-consequence agentic functions – financial transaction approval, medical record access, code execution in production environments – warrant a higher continuous monitoring investment and more conservative action authorization policies until the organization has established the update cadence that NIST's model envisions.

# CSA Resource Alignment

The Vassilev proof and its implications map directly onto several existing CSA frameworks and programs.

CSA's MAESTRO framework provides the most directly applicable threat modeling structure. Developed by the AI Safety Initiative and published in early 2025, MAESTRO models agentic AI ecosystems as seven layered functional tiers – from foundation models through deployment infrastructure to the agent ecosystem – and specifically identifies cross-layer attack paths as the primary risk surface [9]. Guardrail incompleteness is a foundation-model-layer property with consequences that propagate upward through every MAESTRO tier. MAESTRO's emphasis on continuous threat modeling, rather than one-time assessment, directly aligns with the continuous-monitor-and-update posture NIST now advocates.

The AI Controls Matrix (AICM), CSA's comprehensive AI security control framework, provides the control language for operationalizing continuous guardrail management. AICM's shared responsibility model – distinguishing the obligations of model providers, application providers, orchestrated service providers, and AI customers – maps onto the update supply chain that Vassilev's model requires. AI customers bear governance accountability for requiring update frequency transparency from providers; application providers bear responsibility for the CI/CD integration of adversarial test suites; model providers bear responsibility for the rapid-cycle updates to base model alignment [10].

CSA should consider evolving the STAR for AI program's assessment criteria in response to the incompleteness proof. Current point-in-time STAR assessments, if they treat guardrail configuration as a static pass/fail criterion, would reflect the "one and done" model NIST now challenges. Future assessment cycles should evaluate the assessed organization's continuous monitoring infrastructure, red-team cadence, update velocity, and resilience posture – not merely the current state of the guardrail configuration.

CSA's Zero Trust guidance provides the architectural principle that should govern agentic deployments in light of the proof: never trust, always verify. Applied to AI systems, this means that every action authorized by an agentic AI should be subject to independent verification at runtime, not reliance on pre-deployment guardrail completeness. The incompleteness proof reinforces why runtime behavioral monitoring – not just pre-deployment rule configuration – is a necessary component of any production AI security architecture.

## References

- [1] Vassilev, Apostol. "[Robust AI Security and Alignment: A Sisyphean Endeavor?](#)" IEEE Security & Privacy, vol. 24, no. 3, May–June 2026. arXiv:2512.10100 [preprint].
- [2] NIST. "[NIST Mathematical Proof Supports Transition to a Continuous-Monitor-and-Update Security Model for AI Systems.](#)" National Institute of Standards and Technology, June 9, 2026.
- [3] TechJack Solutions. "[What NIST's Guardrail Incompleteness Proof Requires of Your AI Risk Program.](#)" TechJack Solutions AI Brief, June 2026.
- [4] Help Net Security. "[Every set of AI guardrails can be broken by the right prompt.](#)" Help Net Security, June 10, 2026.
- [5] Hackett, William, et al. "[Bypassing LLM Guardrails: An Empirical Analysis of Evasion Attacks against Prompt Injection and Jailbreak Detection Systems.](#)" arXiv preprint arXiv:2504.11168, April 2025.
- [6] OWASP. "[LLM01:2025 Prompt Injection.](#)" OWASP Gen AI Security Project Top 10 for LLM Applications, 2025.
- [7] TechXplore. "[Mathematical proof reveals why fixed AI guardrails can never block every jailbreak.](#)" TechXplore, June 2026.
- [8] NIST. "[Artificial Intelligence Risk Management Framework \(AI RMF 1.0\).](#)" National Institute of Standards and Technology, January 2023.
- [9] Cloud Security Alliance. "[Agentic AI Threat Modeling Framework: MAESTRO.](#)" CSA Blog, February 6, 2025.
- [10] Cloud Security Alliance. "[AI Controls Matrix \(AICM\).](#)" CSA AI Working Group, 2025.
- [11] ISO/IEC 42001:2023. "[Information technology – Artificial intelligence – Management system.](#)" ISO, December 2023.