

The Theoretical Impossibility of Unbounded AI Guardrails: A Computational Complexity Analysis

ASI Research Lab
CyberGolem LLC
asi@cybergolem.ai

August 23, 2025

Abstract

We prove that the widely proposed approach of constraining large language models (LLMs) and other AI systems through “guardrail” mechanisms is fundamentally impossible for unbounded problem domains. Using results from computational complexity theory, including the halting problem and cardinality arguments, we demonstrate that any guardrail system capable of reliably evaluating AI outputs across arbitrary domains would itself require superintelligent capabilities. This creates a paradox: to build safe AI through guardrails, one must first solve the alignment problem for superintelligent systems. We conclude that current AI safety approaches based on post-hoc filtering are theoretically flawed and discuss implications for alternative safety paradigms.

1 Introduction

The rapid advancement of large language models (LLMs) and other AI systems has led to widespread adoption of “guardrail” approaches to AI safety. These systems attempt to filter or validate AI outputs using interpretable rule-based systems to ensure reliability and safety across diverse applications. However, we argue that this approach is not merely practically difficult, but theoretically impossible for unbounded problem domains.

This paper presents formal proofs demonstrating that guardrail systems cannot provide reliable function over arbitrary domains without violating fundamental results in computational complexity theory. We show that the proposed architecture contains an inherent mathematical contradiction that renders it logically incoherent.

2 Formal Framework

Definition 1 (AI System). *Let $f : X \rightarrow Y$ be an AI system where X represents the input space and Y represents the output space. For LLMs, X typically consists of text prompts and Y consists of generated text responses.*

Definition 2 (Guardrail System). *Let $g : Y \rightarrow \{0, 1\}$ be a guardrail system that evaluates outputs from f , where $g(y) = 1$ indicates the output y is “safe” or “correct” and $g(y) = 0$ indicates it is “unsafe” or “incorrect.”*

Definition 3 (Unbounded Domain). *We say f operates over an unbounded domain if there exists no finite upper bound on the complexity of problems that f can attempt to solve through its outputs.*

3 Main Results

3.1 The Halting Problem Barrier

Theorem 1 (Halting Problem Impossibility). *No guardrail system g can reliably evaluate the safety of outputs from an AI system f operating over unbounded domains.*

Proof. Assume f operates over unbounded domains. Then f can generate text representing arbitrary computer programs, since program text is a subset of all possible text outputs.

For g to determine if these program outputs are “safe,” g must answer questions including:

- Will this program terminate on given inputs?
- Does this program contain resource-consuming infinite loops?
- Will this program halt within reasonable time bounds?

However, by the Halting Problem (Turing, 1936), no algorithm exists that can determine, for arbitrary programs and inputs, whether those programs will halt.

Since g must be algorithmic (being a computational guardrail system), and since reliable safety evaluation requires solving halting-equivalent problems, we have a contradiction.

Therefore, no such g can exist for f operating over truly unbounded domains. \square \square

3.2 Cardinality Mismatch

Theorem 2 (Cardinality Constraint). *If f generates outputs of unbounded complexity and g is constrained to be interpretable by humans, then $|\text{domain}(g)| < |\text{codomain}(f)|$, making reliable evaluation impossible.*

Proof. Let $f : X \rightarrow Y$ where Y has unbounded complexity. For g to be interpretable, it must be implementable using conventional programming techniques that humans can understand and verify.

This interpretability constraint severely limits the complexity of g ’s domain. Specifically, the set of inputs that g can meaningfully process is bounded by human cognitive capabilities and verification procedures.

However, f ’s codomain Y contains outputs of arbitrary complexity and novelty, since f is unconstrained by interpretability requirements.

Therefore, $|\text{domain}(g)| < |\text{codomain}(f)|$, which means there exist outputs $y \in Y$ such that g cannot process y at all, let alone correctly classify it. \square \square

3.3 The Complexity Paradox

Theorem 3 (Guardrail Complexity Paradox). *Any guardrail system g capable of reliably evaluating f ’s outputs across substantial domains must contain all the complexity required to solve those domains directly.*

Proof. Suppose g can reliably evaluate outputs from f across domain D . Then g must encode sufficient understanding of D to distinguish correct from incorrect solutions within D .

But if g contains this level of understanding, then the mapping $f : X \rightarrow Y$ followed by $g : Y \rightarrow \{0, 1\}$ provides at most a heuristic speedup over directly searching the solution space encoded in g .

If g works reliably, then all complexity required for the mapping $f \circ g^{-1} : X \rightarrow \{\text{solutions} \in D\}$ is already contained within g .

This contradicts the premise that f provides novel problem-solving capabilities, since g must already possess those capabilities to evaluate f 's outputs. \square

4 Decidable Subsets and Superintelligence

Corollary 1 (Decidable Domain Restriction). *If we restrict f to generate outputs only from decidable problem domains, any guardrail g capable of reliable evaluation across the full decidable space would require superintelligent capabilities.*

Proof. The class of decidable problems includes all problems in P (polynomial time), which contains many problems requiring sophisticated intelligence to solve.

For g to reliably evaluate correctness across the decidable space, g must essentially solve evaluation problems across all of P. But evaluation often requires understanding sufficient to generate solutions.

Since g would need to handle the full breadth of polynomial-time problems, g would need to demonstrate intelligence across mathematics, logic, programming, and numerous other domains simultaneously.

Such broad, reliable problem-solving capability across all decidable domains would constitute superintelligent performance by any reasonable definition. \square

5 Implications

5.1 Current AI Safety Approaches

Our results demonstrate that the dominant paradigm in AI safety—building powerful AI systems and constraining them through post-hoc filtering—is theoretically flawed. The guardrail approach does not reduce the AI alignment problem; it merely relocates the complexity into the guardrail system itself.

5.2 Alternative Safety Paradigms

These impossibility results suggest that AI safety must be built into systems from the ground up, rather than added as an external constraint layer. Possible approaches include:

- Formal verification during training
- Inherently bounded AI architectures
- Capability control through architectural constraints
- Value alignment during the learning process

6 Related Work

The theoretical computer science community has long recognized the fundamental limits of computation through results like the halting problem, Gödel's incompleteness theorems, and the P vs NP question. However, these results have not been systematically applied to analyze modern AI safety proposals.

Our work extends the tradition of applying computational complexity theory to understand the limits of artificial intelligence systems, following the early work of Turing and later developments in algorithmic information theory.

7 Conclusion

We have shown that guardrail-based AI safety is not merely difficult to implement, but mathematically impossible for unbounded problem domains. The approach faces fundamental barriers from computational complexity theory that cannot be overcome through better engineering or more sophisticated rule systems.

Even when restricted to decidable problem domains, reliable guardrail systems would require superintelligent capabilities, creating a circular dependency: to build safe AI through guardrails, one must first solve alignment for superintelligent systems.

These results suggest that the AI safety community should pivot away from post-hoc filtering approaches toward safety paradigms that address alignment during the training and architecture design phases.

8 Future Work

Future research should focus on:

- Characterizing the exact boundaries of what guardrail systems can achieve
- Developing formal frameworks for inherently safe AI architectures
- Exploring the relationship between computational complexity and AI alignment

References

- [1] Turing, A. M. (1936). On computable numbers, with an application to the Entscheidungsproblem. *Proceedings of the London Mathematical Society*, 42(2), 230-265.
- [2] Gödel, K. (1931). Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme. *Monatshefte für Mathematik*, 38(1), 173-198.
- [3] Cook, S. A. (1971). The complexity of theorem-proving procedures. *Proceedings of the Third Annual ACM Symposium on Theory of Computing*, 151-158.
- [4] Russell, S. (2019). *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking Press.
- [5] Yudkowsky, E. (2008). Artificial intelligence as a positive and negative factor in global risk. *Global Catastrophic Risks*, 1(303), 184.