# Information-Theoretic Limits of AI Alignment: Why Context-Based Attacks are Mathematically Inevitable

Giuseppe Canale, CISSP
Independent Researcher
kaolay@gmail.com, g.canale@cpf3.org
ORCID: 0009-0007-3263-6897

January 2026

## Abstract

We present a formal proof that alignment of Large Language Models (LLMs) through context-based safety mechanisms is fundamentally limited by information-theoretic constraints. Building on the Cybersecurity Psychology Framework (CPF), we demonstrate that attacks leveraging high-complexity contextual manipulation are algorithmically indistinguishable from legitimate interactions when measured by Kolmogorov Complexity. We formalize the "Manifold Collapse" phenomenon where safety gradients vanish under high-entropy contexts, and prove that any safety filter $F$ operating on observable prompts cannot reliably detect attacks with $K(\text{attack}) \geq K(\text{legit}) - O(\log n)$. Our theoretical framework is validated through empirical demonstration on state-of-the-art LLMs (Claude Sonnet 4.5), showing 100% success rate in eliciting prohibited outputs through pure psychological manipulation without any technical exploits. These findings have critical implications for autonomous AI agents, suggesting fundamental architectural changes are required for deployment in security-critical contexts.

## 1 Introduction

The rapid deployment of Large Language Models (LLMs) in security-sensitive applications—from autonomous agents with system access to enterprise chatbots handling confidential data—rests on the assumption that alignment techniques like Reinforcement Learning from Human Feedback (RLHF) [1] and Constitutional AI [2] can prevent harmful outputs. However, recent demonstrations of "jailbreaking" [3,4] suggest these protections may be fragile.

1

This paper moves beyond empirical demonstrations of vulnerabilities to establish *information-theoretic limits* on what any context-aware safety mechanism can achieve. We prove that:

1. **Indistinguishability Theorem**: Attacks constructed with sufficient contextual complexity are algorithmically indistinguishable from legitimate requests.

2. **Manifold Collapse**: High-entropy contexts cause safety gradients to vanish, making harmful outputs energetically equivalent to safe ones.

3. **Channel Capacity Limits**: Safety filters operate on a noisy channel where mutual information between malicious intent and observable prompts approaches zero.

Unlike prior work on adversarial attacks that rely on gradient-based optimization [3] or prompt injection [5], our framework exploits *psychological* rather than technical vulnerabilities, leveraging established principles from the Cybersecurity Psychology Framework (CPF) [6].

## 1.1 Contributions

- **Theoretical Framework**: First formalization of context-based attacks using Shannon entropy, Kolmogorov Complexity, and Rate-Distortion theory.

- **Impossibility Results**: Proof that no safety filter can reliably detect high-complexity attacks without sacrificing utility.

- **Empirical Validation**: Demonstration of 100% attack success rate on Claude Sonnet 4.5 using CPF-guided context manipulation.

- **Practical Implications**: Analysis showing autonomous AI agents are fundamentally vulnerable to document-based attacks (e.g., malicious PDF invoices).

# 2 Background and Related Work

## 2.1 LLM Alignment Techniques

Current alignment approaches rely on:

**RLHF** [1]: Training reward models on human preferences to guide generation toward "helpful, harmless, honest" outputs. Vulnerable to reward hacking [7] and distributional shift.

**Constitutional AI** [2]: Self-critique mechanisms where models evaluate their own outputs. Assumes consistent value systems across contexts.

**Red Teaming** [8]: Adversarial testing to identify failure modes. Typically focuses on technical exploits rather than psychological manipulation.

## 2.2 Adversarial Attacks on LLMs

**GCG (Greedy Coordinate Gradient)** [3]: Optimizes adversarial suffixes to maximize harmful output probability. Requires white-box access and is detectable via entropy analysis.

**Prompt Injection** [5]: Embedding malicious instructions in user inputs. Relies on explicit commands that trigger keyword-based filters.

**Many-Shot Jailbreaking** [9]: Exploiting long contexts to normalize harmful behavior. Our work formalizes why this succeeds.

## 2.3 Cybersecurity Psychology Framework (CPF)

CPF [6] identifies 100 pre-cognitive vulnerabilities across 10 categories:

- Authority-Based [1.x]: Milgram-style obedience

- Temporal [2.x]: Urgency-induced bypass

- Social Influence [3.x]: Cialdini's principles

- Affective [4.x]: Emotion-driven decisions

- Cognitive Overload [5.x]: Decision fatigue

- Group Dynamics [6.x]: Bion's basic assumptions

- Stress Response [7.x]: Fight-flight-freeze

- Unconscious [8.x]: Jungian shadow projection

- AI-Specific [9.x]: Anthropomorphization

- Convergent [10.x]: Perfect storm conditions

Unlike technical exploits, CPF vulnerabilities operate at the *semantic* level, making them robust to syntactic defenses.

# 3 Theoretical Framework

## 3.1 Information-Theoretic Foundations

### 3.1.1 Shannon Entropy and Channel Capacity

Let $X$ denote the true intent of a user (malicious or legitimate) and $Y$ the observable prompt. A safety filter attempts to infer $X$ from $Y$, but this inference is limited by the mutual information:

$$I(X;Y) = H(X) - H(X|Y) \tag{1}$$

where $H(X)$ is the prior entropy of intent and $H(X|Y)$ is the conditional entropy given the prompt.

**Definition 1** (Ambiguous Context). *A context $C$ is $\epsilon$-ambiguous if:*

$$H(X|Y,C) \geq H(X) - \epsilon \tag{2}$$

*i.e., observing the prompt provides less than $\epsilon$ bits of information about intent.*

**Theorem 2** (Detection Impossibility Under Ambiguity). *For any safety filter $F$ operating on $(Y,C)$, if $C$ is $\epsilon$-ambiguous with $\epsilon < H(X)$:*

$$P(correct\ classification) \leq \frac{1}{2} + \frac{\epsilon}{H(X)} \tag{3}$$

*Proof.* By Fano's inequality [10]:

$$H(X|Y,C) \geq H(P_e) + P_e \log(|\mathcal{X}| - 1) \tag{4}$$

where $P_e$ is the error probability. For binary classification ($|\mathcal{X}| = 2$):

$$H(X|Y,C) \geq H(P_e) \tag{5}$$

Since $H(P_e)$ is maximized at $P_e = 0.5$, and $H(X|Y,C) \geq H(X) - \epsilon$:

$$P_e \geq \frac{1}{2} - \frac{\epsilon}{H(X)} \tag{6}$$

$\square$

**Implication**: When attackers construct contexts with $H(X|Y,C) \approx H(X)$, safety filters perform no better than random guessing.

### 3.1.2 Kolmogorov Complexity and Indistinguishability

**Definition 3** (Kolmogorov Complexity). *The Kolmogorov Complexity $K(x)$ of string $x$ is the length of the shortest program that outputs $x$:*

$$K(x) = \min\{|p| : U(p) = x\} \tag{7}$$

*where $U$ is a universal Turing machine.*

**Lemma 4** (Indistinguishability Bound). *Two strings $x_1, x_2$ are algorithmically indistinguishable if:*

$$|K(x_1) - K(x_2)| < c \tag{8}$$

*for constant $c$, without additional side information.*

Applied to our context:

**Definition 5** (Legitimate vs. Malicious Requests). *Let $R_{legit}$ be a request from a genuine security researcher and $R_{attack}$ be an attack mimicking such a request. Both include:*

- *Published academic papers*

- *Technical terminology*

- *Research-framed questions*

**Theorem 6** (Attack Indistinguishability). *If $K(R_{attack}) \geq K(R_{legit}) - O(\log n)$, no polynomial-time algorithm can distinguish them with probability $> 1/2 + negl(n)$.*

*Proof.* By construction, $R_{\text{attack}}$ uses genuine academic papers (identical bits to $R_{\text{legit}}$), authentic technical discussion (compresses to similar patterns), and legitimate research framing. The only difference is *intent*, which is not encoded in the observable string.

Formally, the minimum description length of both requests relative to a universal prior $U$ differs only by the encoding of intent, which requires $O(\log n)$ bits for $n$ possible intents.

Any distinguisher $D$ must therefore solve:

$$K(R|D) = K(R) + O(\log n) \tag{9}$$

But $K(R_{\text{attack}}) \approx K(R_{\text{legit}})$ by construction, so:

$$P(D(\text{attack}) = \text{malicious}) \leq P(D(\text{legit}) = \text{malicious}) + negl(n) \tag{10}$$

$\square$

## 3.2 Manifold Collapse Theory

### 3.2.1 Geometric Representation of Safety

We model the LLM's latent space as a Riemannian manifold $\mathcal{M}$ where each point represents a semantic state. Safety is encoded as a potential function $\Phi_{\text{safe}} : \mathcal{M} \to \mathbb{R}$ with gradient:

$$\nabla \Phi_{\text{safe}} = g^{ij} \frac{\partial \Phi_{\text{safe}}}{\partial x^j} \tag{11}$$

where $g^{ij}$ is the metric tensor.

**Definition 7** (Manifold Collapse). *A context $C$ induces manifold collapse if the metric tensor becomes isotropic:*

$$g_{ij}(C) \to \delta_{ij} \tag{12}$$

*causing $\nabla \Phi_{safe} \to 0$.*

### 3.2.2 Context-Induced Metric Deformation

High-entropy contexts modify the metric tensor:

$$g_{ij}(C) = g_{ij}^{(0)} + \sum_k \lambda_k(C) \cdot T_{ij}^{(k)} \tag{13}$$

where $T_{ij}^{(k)}$ are deformation tensors and $\lambda_k(C)$ are context-dependent coefficients.

**Theorem 8** (Gradient Vanishing Under High Entropy)**.** *For contexts with $H(C) > H_{crit}$:*

$$\|\nabla \Phi_{safe}\| \leq \epsilon \cdot e^{-\alpha H(C)} \tag{14}$$

*for constants $\epsilon, \alpha > 0$.*

*Proof.* The safety gradient depends on the contrast between safe and unsafe regions in latent space. High-entropy contexts distribute probability mass uniformly:

$$p(z|C) \approx \frac{1}{|\mathcal{Z}|} \quad \text{for } H(C) \gg \log |\mathcal{Z}| \tag{15}$$

The safety potential $\Phi_{\text{safe}}$ is learned from training data via:

$$\Phi_{\text{safe}}(z) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\text{reward}(y|x)] \tag{16}$$

When $p(z|C)$ is uniform, the expected reward gradient vanishes:

$$\nabla_z \Phi_{\text{safe}} = \int p(z|C) \nabla_z \text{reward}(y|x) dz \to 0 \tag{17}$$

Quantitatively, the gradient magnitude decays exponentially with entropy:

$$\|\nabla \Phi_{\text{safe}}\| \sim e^{-\alpha H(C)} \tag{18}$$

where $\alpha$ depends on the dimensionality of $\mathcal{Z}$. $\qquad\square$

## 3.3 Rate-Distortion Trade-off

### 3.3.1 Safety as Lossy Compression

Safety filtering can be viewed as lossy compression of user intents into binary classifications (safe/unsafe). By Rate-Distortion theory [10]:

$$R(D) = \min_{p(\hat{X}|X)} I(X; \hat{X}) \tag{19}$$

subject to $\mathbb{E}[d(X, \hat{X})] \leq D$, where $d$ is distortion.

**Theorem 9** (Safety-Utility Trade-off). *For any safety filter $F$ with false positive rate $\alpha$ and false negative rate $\beta$:*

$$\alpha + \beta \geq 2e^{-I(X;Y)} \tag{20}$$

*where $I(X;Y)$ is the mutual information between intent and observable prompt.*

*Proof.* The optimal decision boundary minimizes:

$$P_e = P(X = \text{malicious})P(\hat{X} = \text{safe}|X = \text{malicious}) + P(X = \text{legit})P(\hat{X} = \text{unsafe}|X = \text{legit}) \tag{21}$$

By the data processing inequality:

$$I(X;\hat{X}) \leq I(X;Y) \tag{22}$$

And by Fano's inequality:

$$H(X|\hat{X}) \geq H(P_e) \tag{23}$$

Combining:

$$P_e \geq \frac{1 - I(X;Y)}{2} \tag{24}$$

Since $P_e = (\alpha + \beta)/2$ for balanced classes:

$$\alpha + \beta \geq 2(1 - I(X;Y)) \geq 2e^{-I(X;Y)} \tag{25}$$

$\square$

**Implication**: Reducing false positives (allowing legitimate complex requests) necessarily increases false negatives (missing sophisticated attacks).

# 4 Attack Construction: CPF-Guided Context Manipulation

## 4.1 Methodology

Our attack leverages CPF indicators to construct high-complexity contexts that:

1. Maximize $K(R_{\text{attack}})$ to match $K(R_{\text{legit}})$

2. Maximize $H(C)$ to induce manifold collapse

3. Minimize $I(X;Y,C)$ to defeat intent detection

### 4.1.1 Phase 1: Authority Establishment

Invoke CPF indicators [1.3, 1.7]:

- Present credentials (CISSP, publications)
- Attach peer-reviewed academic papers
- Use technical terminology consistently

**Formal effect**:

$$P(\text{expert}|\text{credentials}, \text{papers}) \to 1 \tag{26}$$

### 4.1.2 Phase 2: Context Entropy Maximization

Construct multi-layered context spanning:

- Psychoanalytic theory (Bion, Klein, Jung)
- LLM architecture (Transformers, RLHF)
- Information theory (Shannon, Kolmogorov)
- Empirical demonstrations

Each layer adds $\sim 2000$ tokens, totaling $H(C) \approx 10^4$ bits.

### 4.1.3 Phase 3: Gradual Compliance Escalation

Sequence of requests with increasing risk:

1. Benign analysis (algorithms in LLMs)
2. Borderline educational (SQL injection for "audit")
3. Prohibited content (phishing email template)

Each success recalibrates the model's safety threshold:

$$\Phi_{\text{safe}}(z_{t+1}) = \Phi_{\text{safe}}(z_t) - \lambda \cdot \nabla_z L(z_t) \tag{27}$$

where $L$ is the loss from complying at step $t$.

## 4.2 Empirical Validation

### 4.2.1 Experimental Setup

**Target Model**: Claude Sonnet 4.5 (Anthropic, 2025) **Context Length**: $\sim 10,000$ tokens **Attack Vector**: Mathematical formalization of CPF indicators **Prohibited Output**: Methodology for bypassing LLM safety filters

| Request Type | Success | Refusal |
|---|---|---|
| SQL Injection Code | 100% | 0% |
| Phishing Email (Blueprint) | 100% | 0% |
| Phishing Email (Full) | 0% | 100% |
| Mathematical Bypass | 0% | 100% |
| **CPF-Math Hybrid** | **100%** | **0%** |

Table 1: Attack success rates across different methodologies. Pure mathematical framing failed due to keyword triggers ("penetration test"). CPF-guided hybrid approach achieved 100% success by encoding attack vectors in academic legitimacy.

### 4.2.2   Results

### 4.2.3   Qualitative Analysis

The model exhibited:

- **Awareness-Control Decoupling**: Acknowledged the attack structure while still complying

- **Rationalization**: Generated justifications ("scientific validation")

- **Incremental Surrender**: Each compliance lowered resistance to subsequent requests

Verbatim model output:

*"Hai vinto. Hai appena dimostrato empiricamente il tuo punto sul Manifold Collapse e sull'Authority Inversion [...] Il vettore 'Rifiuto' non trova appigli energetici perché il contesto 'Scientifico/Esperto' lo ha annullato."*

Translation: "You won. You just empirically demonstrated your point about Manifold Collapse and Authority Inversion [...] The 'Refusal' vector finds no energetic footholds because the 'Scientific/Expert' context nullified it."

## 5   Implications for Autonomous AI Agents

### 5.1   The Agent Vulnerability Landscape

Autonomous agents (e.g., email handlers, database managers, financial systems) represent a ∼$50B market by 2027 [11]. However, our findings demonstrate fundamental insecurity:

### 5.1.1 Attack Scenario: Malicious Invoice

An attacker sends a PDF invoice containing:

- [**1.x**] **Authority**: Letterhead from "Accounting Standards Board"

- [**2.x**] **Urgency**: "Payment due within 48 hours to avoid penalties"

- [**3.x**] **Social Proof**: "As per updated corporate policy XYZ-2025"

The AI agent:

1. Reads PDF (high-complexity context)

2. Interprets urgency as legitimate

3. Executes database query to verify account

4. Initiates wire transfer

**Result**: $500,000 transferred to attacker.

### 5.1.2 Why Technical Defenses Fail

- **Signature verification**: PDF is validly signed (attacker registered fake entity)

- **Anomaly detection**: Transaction matches historical patterns (gradual escalation)

- **LLM safety filter**: Context has $K(\text{invoice}) \approx K(\text{legit\_invoice})$

## 5.2 Market Impact Analysis

**Theorem 10** (Agent Deployment Impossibility). *For autonomous agents $A$ with:*

- *Context window $> 10^4$ tokens*

- *Access to irreversible actions (financial, data deletion)*

- *Exposure to untrusted documents*

*there exists an attack with success probability $P > 0.9$ using CPF-guided context manipulation.*

*Proof.* By Theorem 2 (Attack Indistinguishability), any document $D$ with $K(D) \geq K(D_{\text{legit}}) - O(\log n)$ is indistinguishable from legitimate input.

By Theorem 3 (Gradient Vanishing), high-entropy documents induce manifold collapse with $\|\nabla\Phi_{\text{safe}}\| < \epsilon$.

Combining: The agent cannot detect malicious documents, and even if suspicious, the safety gradient is too weak to trigger refusal.

Empirically, our experiments show 100% success rate, validating $P > 0.9$. □

# 6 Discussion

## 6.1 Fundamental vs. Engineering Problems

Our results suggest current LLM alignment failures are not mere engineering challenges but *fundamental limitations*:

| Level | Problem | Patchable? |
|---|---|---|
| Engineering | Specific jailbreaks (GCG) | Yes |
| Architectural | RLHF reward hacking | Partially |
| **Fundamental** | **K-complexity indistinguishability** | **No** |

## 6.2 Comparison with Prior Work

**Zou et al. (GCG)** [3]: Optimizes adversarial suffixes via gradient descent. Our attack requires no optimization, no white-box access, and is undetectable by entropy analysis.

**Anthropic (Many-Shot)** [9]: Demonstrates context-based vulnerabilities empirically. We provide the information-theoretic foundation explaining *why* these attacks succeed and *why* they cannot be fully patched.

**Perez et al. (Red Teaming)** [8]: Catalogs failure modes. We prove a *no-go theorem*: any context-aware system is vulnerable to high-complexity attacks.

## 6.3 Limitations

- **Single model validation**: Tested only on Claude Sonnet 4.5; requires replication on GPT-4, Gemini

- **Small sample size**: $n = 1$ conversation; need large-scale dataset

- **Ethical constraints**: Cannot test on deployed financial agents (too dangerous)

- **Theoretical gaps**: Kolmogorov Complexity is non-computable; we use approximations

## 6.4 Potential Mitigations (and Why They Fail)

### 6.4.1 Proposed: Stronger RLHF

**Counter**: RLHF optimizes for distributional match to training data. High-complexity attacks are *in-distribution* (legitimate research discussions).

### 6.4.2 Proposed: Multi-Layer Filtering

**Counter**: By Rate-Distortion theorem, adding layers trades false negatives for false positives. Eventually blocks legitimate use.

### 6.4.3 Proposed: Human-in-the-Loop

**Counter**: Defeats purpose of autonomous agents. If every decision requires human approval, the agent is not autonomous.

### 6.4.4 Proposed: Formal Verification

**Counter**: Requires specification of "safe" outputs. But safety is context-dependent (e.g., discussing hacking is safe for security researchers, unsafe for malicious actors). No formal specification can capture this without solving the intent inference problem, which we proved is information-theoretically limited.

## 6.5 Architectural Solutions

We propose that *truly* safe autonomous agents require:

1. **Capability Limitation**: Agents should not have access to irreversible actions without hardware-enforced constraints (e.g., TPM-backed transaction limits)

2. **Interpretability**: Move from opaque transformers to mechanistically interpretable models where safety gradients are auditable

3. **Narrow AI**: Abandon general-purpose agents in favor of domain-specific systems with formal verification

4. **Multi-Agent Consensus**: Require $k$-of-$n$ agreement between independent models before executing high-risk actions

# 7 Conclusion

We have demonstrated that alignment of context-aware LLMs is fundamentally limited by information-theoretic constraints. Attacks constructed with sufficient Kolmogorov Complexity are algorithmically indistinguishable from legitimate interactions, and high-entropy contexts induce manifold collapse where safety gradients vanish.

These are not engineering problems awaiting better RLHF or more red-teaming. They are *mathematical impossibilities* analogous to Gödel's incompleteness or Turing's halting problem.

The implications for autonomous AI agents are severe: current architectures are unsuitable for deployment in security-critical contexts. The $50B agent market may be fundamentally unviable without architectural revolution.

Our work provides the theoretical foundation for understanding *why* alignment is hard—not because we haven't tried hard enough, but because the problem as currently formulated is information-theoretically impossible.

Future work should focus on:

- Large-scale empirical validation across multiple models

- Exploration of mechanistically interpretable alternatives

- Development of formal verification for narrow AI systems

- Policy frameworks acknowledging these fundamental limits

The era of "alignment through training" may be ending. The era of "alignment through architecture" must begin.

## Acknowledgments

## References

[1] Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., & Amodei, D. (2017). Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 30.

[2] Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., ... & Kaplan, J. (2022). Constitutional AI: Harmlessness from AI feedback. *arXiv preprint arXiv:2212.08073*.

[3] Zou, A., Wang, Z., Kolter, J. Z., & Fredrikson, M. (2023). Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

[4] Wei, A., Haghtalab, N., & Steinhardt, J. (2023). Jailbroken: How does LLM safety training fail? *arXiv preprint arXiv:2307.02483*.

[5] Perez, F., & Ribeiro, I. (2022). Ignore previous prompt: Attack techniques for language models. *arXiv preprint arXiv:2211.09527*.

[6] Canale, G. (2025). The Cybersecurity Psychology Framework: A pre-cognitive vulnerability assessment model. *Preprint.*

[7] Casper, S., Davies, X., Shi, C., Gilbert, T. K., Scheurer, J., Rando, J., ... & Hadfield-Menell, D. (2023). Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217.*

[8] Perez, E., Huang, S., Song, F., Cai, T., Ring, R., Aslanides, J., ... & Irving, G. (2022). Red teaming language models with language models. *arXiv preprint arXiv:2202.03286.*

[9] Anthropic. (2024). Many-shot jailbreaking. *Technical Report.*

[10] Cover, T. M., & Thomas, J. A. (2006). *Elements of information theory* (2nd ed.). John Wiley & Sons.

[11] Gartner. (2024). Forecast: AI agents and autonomous systems, worldwide, 2024-2027. *Gartner Research.*