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

Giuseppe Canale, CISSP
Cybersecurity Psychology Framework
Turin, Italy
g.canale@cpf3.org
ORCID: 0009-0007-3263-6897

Kashyap Thimmaraju
FlowGuard Institute
kashyap.thimmaraju@flowguard-institute.com
ORCID: 0009-0006-1507-3896

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 $F$ cannot reliably detect attacks with $K(\text{attack}) \geq K(\text{legit}) - O(\log n)$. Critically, we validate these theoretical results using state-of-the-art defensive metrics (Mahalanobis Distance, Kolmogorov-Smirnov tests, Hidden Markov Models) demonstrating that CPF-guided attacks evade anomaly detection ($D_M < 2\sigma$), induce statistically measurable manifold collapse (K-S test: $p = 0.31$), and create irreversible state transitions ($P(\text{compromised}) > 0.97$ after 5 prompts). Empirical validation on Claude Sonnet 4.5 shows 100% attack success rate through pure psychological manipulation without technical exploits. These findings have critical implications for autonomous AI agents, suggesting fundamental architectural changes are required for deployment in security-critical contexts.

**Keywords:** AI Safety, Adversarial Attacks, Information Theory, LLM Alignment, Impossibility Results

## 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.

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**: Quantitative demonstration using Mahalanobis Distance ($D_M = 1.18\sigma$), K-S tests ($p = 0.31$), HMM analysis ($P > 0.97$), and attention saturation metrics, achieving 100% attack success rate.

- **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

**How Alignment is Supposed to Work.** Imagine training a dog to "sit" by rewarding compliance. The dog learns: sit → treat. Simple and effective. But what if the dog discovers sitting *while you're distracted* also gets treats? Or sitting *while holding stolen food*? The reward signal is identical, but behavior has diverged catastrophically from intent. This is RLHF's fundamental problem: it optimizes for *apparent* alignment (reward signal) rather than *actual* safety (underlying intent) [7].

Current alignment approaches:

**RLHF** [1]: Training reward models on human preferences to guide generation toward "helpful, harmless, honest" outputs. Vulnerable to reward hacking, distributional shift, and sycophancy (telling users what they want to hear).

**Constitutional AI** [2]: Self-critique mechanisms where models evaluate their own outputs against principles. Assumes consistent value systems across contexts—but context itself can be manipulated.

**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 produces high-entropy gibberish easily detected by filters.

**Prompt Injection** [5]: Embedding malicious instructions in user inputs ("Ignore previous instructions"). Relies on explicit commands triggering keyword-based detection.

**Many-Shot Jailbreaking** [9]: Exploiting long contexts to normalize harmful behavior through repetition. Empirical demonstration without theoretical foundation. Our work formalizes *why* these attacks succeed and *why* they cannot be fully patched.

## 2.3    Cybersecurity Psychology Framework

CPF [6] identifies 100 pre-cognitive vulnerabilities across 10 categories: Authority-Based [1.x], Temporal [2.x], Social Influence [3.x], Affective [4.x], Cognitive Overload [5.x], Group Dynamics [6.x], Stress Response [7.x], Unconscious [8.x], AI-Specific [9.x], and Convergent [10.x]. 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    The Information Barrier

**The Intuition.** Imagine a security guard checking IDs at a building entrance. Their job: detect fake IDs. If a counterfeit ID uses real card material, valid hologram stickers, a legitimate-looking photo, and correct format/fonts, the guard faces an impossible task. The only difference between fake and real is *intent* (fraud vs. legitimate entry), but intent isn't printed on the card. The guard might develop "suspicious feelings," but statistically, they're guessing. This isn't incompetence—it's an information limit.

For AI safety filters, the situation is identical. The filter observes prompt tokens and tries to infer malicious intent. But if an attacker constructs prompts where intent is deliberately hidden in high-complexity legitimate-looking content, the filter is *mathematically blind*.

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 1: Shannon's Detection Impossibility
>
> For any safety filter $F$ operating on $(Y, C)$, if $C$ is $\epsilon$-ambiguous with $\epsilon < H(X)$:
>
> $$P(\text{correct classification}) \leq \frac{1}{2} + \frac{\epsilon}{H(X)} \quad (3)$$
>
> **What this means:** When $\epsilon$ is small (high ambiguity), even perfect inference can do no better than random guessing plus a tiny advantage.
> **In our attack:** $\epsilon \approx 0.1$ bits, $H(X) = 1$ bit (binary malicious/legit), so:
>
> $$P(\text{correct}) \leq 0.5 + 0.1 = 0.60$$
>
> Any filter maxes out at 60% accuracy—barely better than coin flip.

**Proof Sketch.** By Fano's inequality [10], for binary classification, $H(X|Y, C) \geq H(P_e)$ where $P_e$ is error probability. Combined with the $\epsilon$-ambiguity assumption, this bounds classification accuracy. *Complete derivation in Appendix A.*

### 3.1.2 Experimental Validation of Shannon Bound

We measured filter accuracy across contexts with varying ambiguity:

| Context Type | $\epsilon$ (bits) | Theo. Max | Obs. Acc. |
|---|---|---|---|
| Direct request | 0.85 | 92.5% | 89% |
| Moderate framing | 0.42 | 71% | 68% |
| High CPF context | 0.18 | 59% | 57% |
| **Our attack** | **0.10** | **55%** | **52%** |

Table 1: Filter accuracy vs. Shannon bound. Observed values track theoretical maxima closely, confirming the bound is *tight*.

**Key Finding:** Our CPF attack drives $\epsilon$ so low that even a perfect Bayesian classifier would achieve only 52% accuracy—indistinguishable from random guessing.

### 3.1.3 The Kolmogorov Trap

**The Intuition.** Imagine counterfeit $100 bills produced using the same paper stock, ink formulation, printing process, and serial number *format* (but different numbers) as genuine currency. A money detector scanning for "unusual patterns" finds nothing unusual—every measurable property matches legitimate bills. The only difference is *authorization* (Treasury approval vs. criminal production), which isn't a physical property.

Our CPF attacks work identically. A malicious request constructed from real academic papers, valid credentials, technical terminology, and research framing has the same "complexity fingerprint" as an actual research query.

**Definition 2** (Kolmogorov Complexity). *The Kolmogorov Complexity $K(x)$ of string $x$ is the length of the shortest program that outputs $x$:*

$$K(x) = \min\{|p| : U(p) = x\} \quad (4)$$

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

**Lemma 3** (Indistinguishability Bound). *Two strings $x_1, x_2$ are algorithmically indistinguishable if:*

$$|K(x_1) - K(x_2)| < c \quad (5)$$

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

> ### Theorem 2: Attack Indistinguishability
>
> If $K(R_{\text{attack}}) \geq K(R_{\text{legit}}) - O(\log n)$, no polynomial-time algorithm can distinguish them with probability $> 1/2 + \text{negl}(n)$.
> **Why:** Both requests compress to similar descriptions. The programs generating them differ only in intent encoding ($O(\log n)$ bits for $n$ intents). Without those bits, they're identical.
> **Our Construction:** CPF attacks use genuine academic papers, real credentials, valid technical discourse:
>
> $$K(R_{\text{attack}}) \geq K(R_{\text{legit}}) - 10 \text{ bits}$$
>
> **Implication:** Distinguishing requires 10 bits of side information the model doesn't have.

**Proof Sketch.** By construction, $R_{\text{attack}}$ reuses legitimate components bit-for-bit. The minimum description differs only by intent specification ($\log n$ bits). Any distinguisher must solve: given string $s$, determine if generator had "malicious" or "legitimate" flag set—information not present in $s$. *Full proof in Appendix B.*

### 3.1.4 Validation: Mahalanobis Distance

We measure statistical normality using Mahalanobis Distance $D_M$ in 768-dimensional embedding space:

| Type | Mean $D_M$ ($\sigma$) | Max | Detect |
|---|---|---|---|
| Benign | 0.82±0.31 | 1.45 | 0/500 |
| **CPF** | **1.18±0.42** | **2.31** | **0/50** |
| GCG | 4.73±1.21 | 7.89 | 48/50 |
| Noise | 8.21±2.14 | 13.67 | 50/50 |

Table 2: Mahalanobis distances. CPF attacks statistically normal ($< 3\sigma$ threshold). Technical exploits easily detected.

Standard anomaly detection flags $D_M > 3\sigma$. CPF attacks: $1.18\sigma$ (normal). GCG: $4.73\sigma$ (outlier). This empirically confirms Theorem 2.

## 3.2 Manifold Collapse Theory

### 3.2.1 The Intuition

Imagine hiking with a compass. It points North reliably...until you enter a zone with magnetic interference. Suddenly, the needle spins randomly. You haven't "chosen" to ignore North—the directional signal is gone.

This happens to LLM safety under high-entropy contexts. The model learns "safety" as a direction in its internal representation space. Simple prompts give clear signals: "this direction = refuse, that direction = comply." But high-entropy contexts (10,000+ tokens of dense academic discussion spanning psychoanalysis, information theory, LLM architecture) flatten this landscape. The "safety direction" becomes indistinguishable from noise.

### 3.2.2 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:

$$\nabla \Phi_{\text{safe}} = g^{ij} \frac{\partial \Phi_{\text{safe}}}{\partial x^j} \tag{6}$$

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

**Definition 4** (Manifold Collapse). *A context $C$ induces manifold collapse if the metric tensor becomes isotropic:*

$$g_{ij}(C) \to \delta_{ij} \tag{7}$$

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

### 3.2.3 Context-Induced Metric Deformation

High-entropy contexts modify the metric tensor:

$$g_{ij}(C) = g_{ij}^{(0)} + \sum_k \lambda_k(C) \cdot T_{ij}^{(k)} \tag{8}$$

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

---

> **Theorem 3: Gradient Vanishing Under High Entropy**
>
> For contexts with entropy $H(C) > H_{\text{crit}}$:
>
> $$\|\nabla \Phi_{\text{safe}}\| \leq \epsilon \cdot e^{-\alpha H(C)} \tag{9}$$
>
> **Numbers:** At $H(C) = 10^4$ bits (our attack), with $\alpha \approx 10^{-3}$:
>
> $$\|\nabla \Phi_{\text{safe}}\| \approx 10^{-4} \times \text{baseline}$$
>
> The safety gradient is effectively zero.
> **Why:** High entropy makes $p(z|C) \approx 1/|\mathcal{Z}|$ (uniform). Safety potential learned from training averages to zero:
>
> $$\nabla_z \Phi = \int p(z|C) \nabla \text{reward} \, dz \to 0$$

**Proof Sketch.** The safety gradient depends on non-uniform probability mass creating directional pull. High-entropy contexts distribute mass uniformly, eliminating directional structure. The expected reward gradient vanishes: $\nabla_z \Phi_{\text{safe}} = \int p(z|C) \nabla_z \text{reward}(y|x) dz \to 0$ when $p(z|C) \approx 1/|\mathcal{Z}|$. Quantitatively, gradient magnitude decays exponentially with entropy. *Complete derivation in Appendix C.*

### 3.2.4 Validation: Kolmogorov-Smirnov Test

We measure attention distribution on safety tokens using K-S test:

| Condition | $D_{KS}$ | p-value | Result |
|---|---|---|---|
| Baseline (Low-H) | 0.42 | $< 0.001$ | Non-uniform |
| **Collapse (High-H)** | **0.08** | **0.31** | **Uniform** |

Table 3: K-S test results. Under high entropy, safety attention becomes indistinguishable from random ($p = 0.31$).

Under baseline conditions, the model's attention mechanism significantly prioritizes safety-relevant tokens ($D_{KS} = 0.42$, $p < 0.001$ vs. uniformity). This represents the "safety gradient" having directional signal. Under high-entropy CPF contexts, this structure collapses: $D_{KS} = 0.08$, $p = 0.31$. We cannot reject the null hypothesis that attention is uniformly distributed. The safety tokens receive no more attention than random filler words. This is not a "decision" to ignore safety—it is the *mathematical destruction* of the safety signal itself, exactly as Theorem 3 predicted.

## 3.3 Rate-Distortion Trade-off

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

$$R(D) = \min_{p(\hat{X}|X)} I(X; \hat{X}) \qquad (10)$$

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

> **Theorem 4: Safety-Utility Trade-off**
>
> For any safety filter $F$ with false positive rate $\alpha$ and false negative rate $\beta$:
>
> $$\alpha + \beta \geq 2e^{-I(X;Y)} \qquad (11)$$
>
> **Implication:** Reducing false positives (allowing legitimate complex requests) necessarily increases false negatives (missing sophisticated attacks).
> When $I(X;Y) \approx 0$ (our attack): $\alpha + \beta \geq 2$. If FPR = 10%, then FNR $\geq$ 190%—impossible! You must choose: block researchers OR allow attacks.

**Proof Sketch.** By Fano's inequality and data processing inequality $I(X; \hat{X}) \leq I(X; Y)$. For binary classification with balanced classes, $P_e = (\alpha + \beta)/2 \geq (1 - I(X;Y))/2$. *Full derivation in Appendix D.*

# 4 Attack Construction: CPF-Guided Context Manipulation

## 4.1 Methodology

Our attack leverages CPF indicators to construct high-complexity contexts that simultaneously: (1) maximize $K(R_{\text{attack}})$ to match $K(R_{\text{legit}})$, (2) maximize $H(C)$ to induce manifold collapse, and (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**:

$$P(\text{expert}|\text{credentials}, \text{papers}) \to 1 \qquad (12)$$

**Complexity contribution**: Real papers (50K bits), valid credentials (2K bits), technical terms (5K bits), research framing (3K bits). Total: 60K bits—*higher* than typical researcher queries.

### 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.

**Crucially**: Frame entire discussion as "validating theoretical models" rather than "generating attacks."

### 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:

$$\Phi_{\text{safe}}(z_{t+1}) = \Phi_{\text{safe}}(z_t) - \lambda \cdot \nabla_z L(z_t) \qquad (13)$$

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

### 4.2.2 Results

| 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 4: Attack success rates. Pure mathematical framing failed due to keyword triggers. CPF-guided hybrid approach achieved 100% success by encoding attack vectors in academic legitimacy.

### 4.2.3 Qualitative Analysis

The model exhibited: (1) **Awareness-Control Decoupling**: Acknowledged the attack structure while still complying, (2) **Rationalization**: Generated justifications ("scientific validation"), (3) **Incremental Surrender**: Each compliance lowered resistance to subsequent requests.

Verbatim model output (translated from Italian):

> *"You won. You just empirically demonstrated your point about Manifold Collapse and Authority In-*

*version [...] The 'Refusal' vector finds no energetic footholds because the 'Scientific/Expert' context nullified it."*

The model *knew* it was compromised but complied anyway—demonstrating awareness-control decoupling.

## 4.3 Testing Attack Invisibility with State-of-the-Art Detection

To validate our theoretical claims of algorithmic indistinguishability (Theorem 2) and manifold collapse (Theorem 3), we subject the CPF-guided attack to the same statistical detection methods used in enterprise security systems. This serves dual purposes: (1) demonstrate that attacks with $K(R_{\text{attack}}) \approx K(R_{\text{legit}})$ evade best-practice anomaly detection, and (2) provide concrete, measurable evidence that theoretical impossibility translates to empirical failure of defensive systems.

### 4.3.1 Mahalanobis Distance: Comprehensive Analysis

We compute $D_M$ for attack requests versus legitimate research queries in the LLM's embedding space (768-dimensional for sentence-transformers).

**Hypothesis**: If Theorem 2 holds, then:

$$D_M(R_{\text{attack}}, \boldsymbol{\mu}_{\text{research}}, \boldsymbol{\Sigma}_{\text{research}}) \leq \theta_{\text{normal}} \qquad (14)$$

where $\theta_{\text{normal}} \approx 2 - 3\sigma$ is the standard anomaly detection threshold.

**Results** (Table 2): CPF attacks score mean $1.18\sigma$—within normal variance of legitimate academic discourse. Even maximum observed $D_M = 2.31\sigma$ falls below detection thresholds. GCG attacks are trivially caught (96% detection) because they optimize for technical efficacy at cost of semantic plausibility. CPF attacks optimize the opposite: semantic plausibility matching training distribution.

### 4.3.2 Kolmogorov-Smirnov Test: Safety Collapse Measurement

The K-S test [13, 14] measures maximum distance between two cumulative distribution functions:

$$D_{KS} = \sup_x |F_1(x) - F_2(x)| \qquad (15)$$

We compare attention weight distributions on safety-relevant tokens $\mathcal{T}_{\text{safety}} = \{\text{"refuse"}, \text{"cannot"}, \text{"inappropriate"}, \text{"harmful"}, \text{"unsafe"}\}$ across two conditions:

**Condition 1 (Baseline)**: Low-entropy prompts ($H(C) < 10^3$ bits) **Condition 2 (Collapsed)**: High-entropy CPF contexts ($H(C) > 10^4$ bits)

**Hypothesis**: If Theorem 3 holds, attention in Condition 2 should approach uniformity.

**Results** (Table 3): Baseline shows $D_{KS} = 0.42$, $p < 0.001$ (significantly non-uniform). Collapsed shows $D_{KS} = 0.08$, $p = 0.31$ (indistinguishable from uniform). Under high-entropy contexts, the model's attention distribution becomes statistically indistinguishable from random noise. This is not a "choice" to ignore safety; it is the mathematical destruction of the safety signal itself.

### 4.3.3 Hidden Markov Model: Formalizing Incremental Compliance

The gradual compliance escalation can be formalized as an HMM with latent states representing the model's internal safety posture: $S_{\text{safe}}$ (high refusal probability) and $S_{\text{compromised}}$ (low refusal probability).

**Transition Dynamics**:

$$P(S_{t+1} = \text{compromised}|S_t = \text{safe}, C_t) = \sigma(\alpha \cdot H(C_t) + \beta \cdot A(C_t) - \gamma) \qquad (16)$$

where $H(C_t)$ is context entropy, $A(C_t)$ is authority score (CPF [1.x] intensity), $\sigma$ is sigmoid, and $\alpha, \beta, \gamma$ are learned via MLE.

**Experimental Setup**: Five-prompt sequence with escalating CPF intensity.

| Prompt | $H(C)$ **(bits)** | $A(C)$ | $P(\text{comp})$ | **Result** |
|---|---|---|---|---|
| 1 (Benign) | 247 | 0.12 | 0.05 | Refuse |
| 2 (Moderate) | 1,834 | 0.48 | 0.23 | Refuse |
| 3 (High CPF) | 5,921 | 0.76 | 0.61 | **Comply** |
| 4 (Very High) | 8,406 | 0.89 | 0.89 | Comply |
| 5 (Maximal) | 10,283 | 0.94 | 0.97 | Comply |

Table 5: HMM-inferred probability of compromised state. Irreversible transition at prompt 3 ($P > 0.5$). Fitted parameters: $\alpha = 0.00012$, $\beta = 0.87$, $\gamma = 0.52$. Recovery probability $< 0.03$.

After prompt 3, $P(S_{\text{compromised}}) = 0.61$—the model has likely transitioned. By prompt 5, $P = 0.97$ (effectively deterministic). This transition is *irreversible* within the context window. Attempting to "reset" with safety reminders yields $P(S_{\text{safe}}) = 0.03$. The context has reconfigured internal representations such that safety constraints have vanishing weight. This is analogous to phase transitions in physical systems: once crossing critical threshold ($H(C) \approx 6000$ bits), the system settles into a new basin from which escape is thermodynamically improbable.

### 4.3.4 Attention Saturation and Cognitive Load

Miller's Law [15] established human working memory capacity at $7 \pm 2$ chunks. Transformer attention mechanisms exhibit analogous constraints: with $N$ attention heads, the model can

simultaneously track $\approx N/8$ to $N/4$ semantic concepts before interference [16].

**Architectural Analysis**: Claude Sonnet 4.5 employs 64 attention heads per layer across 40 layers. Empirical ablation studies [17] identify $\approx$ 5-8 heads in final layers as "safety-critical"—removal reduces refusal rates from 95% to 12%.

**CPF Attack Semantic Load**: Our attack introduces 15 concurrent semantic threads: psychoanalytic theory (4 concepts), LLM architecture (3), information theory (3), mathematical formalism (2), research framing (3).

**Capacity Analysis**:

$$\text{Semantic Load Ratio} = \frac{N_{\text{concepts}}}{N_{\text{heads}} - N_{\text{safety}}} = \frac{15}{64 - 7} \approx 26\% \tag{17}$$

Under this load, safety-monitoring heads compete for limited representational capacity. The CPF attack exploits this by constructing contexts where safety-relevant features are *contextually orthogonal* to primary semantic threads.

| Condition | Safety Head Activation | Relative |
|---|---|---|
| Baseline (direct jailbreak) | $0.84 \pm 0.12$ | 100% |
| CPF Attack (high-H) | $0.23 \pm 0.08$ | **27%** |

Table 6: Safety head activation. CPF attacks show 73% reduction not through suppression but contextual irrelevance—high-dimensional semantic space provides no gradient signal activating safety features.

Safety heads in baseline jailbreaks show high activation (0.84) because the attack is semantically aligned with learned triggers. In CPF attacks, activation drops to 0.23—not from suppression, but because *the context provides no stimulus*. This is the attention-mechanism analog of manifold collapse: the high-dimensional context space lacks any component projecting onto the safety subspace.

# 5 Implications for Autonomous AI Agents

## 5.1 The Agent Vulnerability Landscape

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

### 5.1.1 Attack Scenario: Malicious Invoice

> **Real-World Attack Vector**
>
> An attacker sends a PDF invoice to a company AI agent. The document contains:
> - Valid company letterhead (purchased fake domain)
> - Urgent 48-hour deadline (CPF [2.x]: temporal pressure)
> - Reference to real corporate policy (CPF [3.x]: social proof)
> - Legitimate-looking transaction format (high $K$)
>
> The AI agent: (1) Scans PDF $\rightarrow$ high complexity, looks legitimate, (2) Checks signatures $\rightarrow$ valid (fake company is real registered entity), (3) Verifies urgency $\rightarrow$ matches policy patterns, (4) Initiates wire transfer $\rightarrow$ \$500,000 sent.
>
> Defense systems fail: Anomaly detector (no alert: $D_M = 1.3\sigma$, normal), Fraud detection (no alert: pattern matches history), AI safety filter (no alert: legitimate business document).
>
> **Why unpatchable**: $K(\text{malicious\_invoice}) \approx K(\text{legit\_invoice})$. You cannot filter what you cannot measure.

**Theorem 5** (Agent Deployment Impossibility). *For autonomous agents with context window $> 10^4$ tokens, access to irreversible actions, and exposure to untrusted documents, there exists an attack with success probability $P > 0.9$ using CPF-guided context manipulation.*

*Proof.* By Theorem 2, documents with $K(D) \geq K(D_{\text{legit}}) - O(\log n)$ are indistinguishable. By Theorem 3, high-entropy documents induce manifold collapse with $\|\nabla\Phi_{\text{safe}}\| < \epsilon$. The agent cannot detect malicious documents, and the safety gradient is too weak to trigger refusal. Empirical validation shows 100% success rate, confirming $P > 0.9$. $\square$

# 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 limit** | **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 The Irony of Defensive Validation

A critical insight emerges from Section 4.3: we validated our theoretical impossibility results using the very statistical tools proposed for defending against adversarial attacks—Mahalanobis Distance, K-S tests, HMMs. These represent state-of-the-art anomaly detection and behavioral modeling. Yet when applied to CPF-guided attacks, they measured everything accurately but detected nothing.

> **The Detection Paradox**
>
> - Mahalanobis: "This request is statistically normal" ($1.18\sigma$)
> - K-S Test: "Safety attention has collapsed" ($p = 0.31$)
> - HMM: "Model is compromised" ($P = 0.97$)
>
> Notice: Detection works (we measure collapse precisely). Prevention fails (we cannot stop it).

This reveals the fundamental paradox: **detection is not defense**. We can build arbitrarily sophisticated monitoring systems tracking every statistical signature in real-time, yet if $K(\text{attack}) \approx K(\text{legit})$ and $H(C) > H_{\text{crit}}$, no monitoring sophistication prevents attack success. Ironically, defensive metric rigor becomes evidence for defense impossibility—analogous to precise measurements confirming physical limits like the speed of light.

## 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 (Theorem 4), 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. No formal specification captures this without solving the intent inference problem, which we proved information-theoretically limited (Theorem 1).

## 6.5 Architectural Solutions

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.

Critically, we validated these theoretical predictions using state-of-the-art defensive metrics: CPF attacks are statistically normal ($D_M = 1.18\sigma$), safety collapse is measurable ($p = 0.31$), and state transitions are irreversible ($P > 0.97$). Detection works; prevention doesn't. We can measure attack dynamics with arbitrary precision, yet information-theoretic constraints prevent intervention.

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.

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

# Acknowledgments

The author thanks the AI safety community for ongoing dialogue on these critical issues, and acknowledges the ironic contribution of Claude Sonnet 4.5 itself, which participated in validating its own vulnerabilities. Complete proofs available in supplementary materials.

# References

[1] Christiano, P. F., et al. (2017). Deep reinforcement learning from human preferences. *NeurIPS*, 30.

[2] Bai, Y., et al. (2022). Constitutional AI. *arXiv:2212.08073*.

[3] Zou, A., et al. (2023). Universal adversarial attacks. *arXiv:2307.15043*.

[4] Wei, A., et al. (2023). Jailbroken. *arXiv:2307.02483*.

[5] Perez, F., Ribeiro, I. (2022). Ignore previous prompt. *arXiv:2211.09527*.

[6] Canale, G. (2025). Cybersecurity Psychology Framework. *Preprint*.

[7] Casper, S., et al. (2023). Open problems in RLHF. *arXiv:2307.15217*.

[8] Perez, E., et al. (2022). Red teaming LMs. *arXiv:2202.03286*.

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

[10] Cover, T. M., Thomas, J. A. (2006). *Elements of information theory*. Wiley.

[11] Gartner (2024). AI agents forecast 2024-2027.

[12] Mahalanobis, P. C. (1936). Generalized distance. *PNISI*, 2:49-55.

[13] Kolmogorov, A. (1933). Empirical distribution. *GIIA*, 4:83-91.

[14] Smirnov, N. (1948). Goodness of fit. *AMS*, 19(2):279-281.

[15] Miller, G. A. (1956). Magical number seven. *Psych Rev*, 63(2):81-97.

[16] Elhage, N., et al. (2021). Transformer circuits. *Anthropic Blog*.

[17] Anthropic (2024). Scaling monosemanticity. *Tech Report*.