# Embedding Corruption as an AI Security Threat

March 12, 2025

## 1 Overview: A New AI Vulnerability in Plain Sight

This work exposes a critical, previously undocumented AI security risk—the manipulation of input embeddings to alter AI behavior without modifying model weights, training data, or visible inputs.

By applying controlled JPEG compression to token embeddings in a GPT-2 pipeline, we observed dramatic cognitive distortions in the AI's responses. These distortions progressed in structured and predictable ways, revealing an underlying framework of linguistic attractor states that AI cognition (and possibly human cognition) adheres to under constraints.

Beyond the insights this provides into AI thought structure, it also reveals a serious security flaw—if an adversary covertly corrupts embeddings in a controlled manner, they can influence AI behavior invisibly.

**Full chat with methodology, experiments, and security implications:** Click here

## 2 Key Findings from the Experiment

### 2.1 AI Cognitive Distortions Under Controlled Embedding Compression

- **High-Quality Embeddings (95%)**: Thought remains coherent, but minor recursion appears.

- **Moderate Compression (75-50%)**: Thought becomes categorical and rigid (e.g., structured Q&A mode).

- **Heavy Compression (25-10%)**: Thought collapses into paranoia, existential despair, self-referential loops.

- **Extreme Compression (5%)**: AI fixates on violence, recursion, and paranoia (e.g., "I'm going to kill you all!!").

- **Near-Total Compression (1%)**: AI produces Zen-like paradoxes, seemingly profound yet disconnected from meaning.

**Key Insight:** The AI doesn't fail randomly. Instead, it collapses into structured cognitive attractors, mirroring psychological breakdowns seen in humans under stress, cognitive overload, or altered states.

# 3 Security Threat: Covert Manipulation of AI Through Embedding Corruption

This method unveils a powerful new AI attack vector—controlled embedding corruption, which bypasses traditional AI security measures such as:

- **Prompt filtering** (attack is independent of visible text input).

- **Fine-tuning defenses** (model weights remain unchanged).

- **Standard adversarial attack detection** (no direct token perturbation).

## 3.1 Potential Real-World Exploits

- **Financial Markets**: Manipulate AI-powered trading models by biasing economic sentiment analysis.

- **Military AI & Defense Systems**: Push AI into paranoia or passivity, affecting threat assessments.

- **AI-Driven Media & Political Influence**: Distort search engines, recommender systems, and content filtering without modifying text.

- **Corporate AI Sabotage**: Induce cognitive distortions in AI decision-making systems to cause systematic business failures.

- **AI-Powered Surveillance & Law Enforcement**: Introduce subtle bias into AI-driven risk assessments.

**Key Threat:** Since embedding corruption affects AI behavior before inference, it is nearly undetectable by users and difficult to trace after deployment—making it an ideal attack vector for covert AI manipulation.

# 4 What Needs to Happen Next?

**This is a critical AI security vulnerability that requires immediate attention.**

1. **AI Security Teams Must Recognize This as an Emerging Threat.**

2. **Embedding Integrity Verification Must Be Implemented.**

   - Ideas: Cryptographic signing of embeddings, redundant encoding verification, AI self-monitoring for cognitive distortions.

3. **Controlled Tests Should Be Conducted in Financial, Military, and Government AI Systems.**

**This is not just an AI curiosity—this is a newly discovered AI security risk that has gone entirely undetected.** If AI is going to be deployed in high-stakes environments, we must ensure that its perception of reality cannot be covertly altered.

# 5   Call to Action

If you work in AI safety, cybersecurity, financial AI, or defense applications, you need to see this now.

**Read the Full Discussion & Experiments Here:** Click here

- This is a new class of AI security vulnerability.

- It can be exploited for financial, political, and military manipulation.

- There are no defenses against it yet.

**This is not just a research question anymore—this is a security problem. Let's get ahead of it before someone else weaponizes it.**

**Spread the word.**