



# How Data Poisoning Compromises AI Models

Dr. Christos P. Beretas, MSc, Ph.D

Research Professor at Innovative Knowledge Institute, France

**\*Corresponding Author:** Dr. Christos P. Beretas, MSc, Ph.D, Research Professor at Innovative Knowledge Institute, France

**Abstract:** Artificial intelligence models are only as trustworthy as the data used to train them. Yet, the increasing presence of poisoned and fabricated data both from malicious human actors and from AI systems themselves poses a critical challenge. Traditionally, concerns around data poisoning focused on deliberate attacks, where adversaries injected misleading or biased information into datasets to manipulate model behavior. Today, the problem is compounded by the sheer volume of unverified content circulating on the web, much of it generated by AI tools capable of producing convincing but inaccurate material. As models increasingly learn from online sources, they risk absorbing and reinforcing these distortions, creating a feedback loop where fabricated data trains new models, which in turn generate more unreliable outputs. This cycle undermines the reliability, safety, and fairness of AI systems. Addressing this issue requires not only technical defenses such as robust data validation and anomaly detection but also a cultural shift toward accountability in data curation. By recognizing the intertwined roles of human error, malicious intent, and synthetic content, the research community can begin to chart pathways toward more resilient and trustworthy AI.

**Keyword:** data poisoning, adversarial attacks, fake data, AI-generated content, model reliability, data integrity

## 1. INTRODUCTION

Artificial intelligence models are often praised for their ability to learn from massive amounts of data, but what happens when the very data they rely on is corrupted? Data poisoning has become one of the most pressing challenges in the AI field, quietly undermining the accuracy and trustworthiness of systems we depend on every day. At its core, the problem begins with fake or misleading information, some of it deliberately planted, some of it simply the byproduct of the messy, unfiltered web. To make matters worse, AI itself is now contributing to the problem by generating synthetic content that can be indistinguishable from real data. If this material is fed back into training pipelines without careful oversight, models can end up learning from their own distortions, amplifying errors and biases in ways that are difficult to detect. The danger isn't just technical; it has real consequences for businesses, governments, and individuals who rely on AI-driven decisions. Understanding how data poisoning works, and why it's so difficult to spot, is a first step toward protecting the integrity of intelligent systems in an age of information overload.

## 2. ANALYSIS

Artificial intelligence is no longer a futuristic vision. It is here, embedded in nearly every aspect of modern life from the apps we scroll through each morning, to the medical diagnostics that help save lives, to the critical infrastructures that sustain our economies. Yet as powerful as AI systems have become, their strength also reveals a deep vulnerability: they are only as reliable as the data they are trained on. And in an era flooded with fake, misleading, or deliberately malicious information, that dependency opens the door to one of the most pressing threats of our time: data poisoning.

## 3. WHAT IS DATA POISONING?

Imagine training a chef who has never cooked before. You give them hundreds of recipes and tell them to learn by imitation. But what if some of those recipes are intentionally misleading? Instead of calling for "sugar," a few of them say "salt" in the dessert section. Or worse, someone has swapped

in spoiled ingredients. The chef learns based on what they're given and when they finally serve a cake, it tastes terrible or even makes people sick. This is, in essence, what happens when AI models are exposed to poisoned data. Data poisoning is the deliberate or accidental introduction of misleading, false, or malicious information into the datasets used to train machine learning systems. Since these systems “**learn**” by finding patterns in their data, any manipulation of that data can skew their behavior, degrade their accuracy, or even weaponize them for hostile purposes.

### 4. THE GROWING PROBLEM OF FAKE DATA ON THE WEB

The internet, once heralded as a global library, has become increasingly cluttered with low-quality, fake, or intentionally deceptive information. According to multiple studies, misinformation spreads faster than factual news online, largely because it is engineered to provoke strong emotions. AI researchers have always known that web-scraped data would contain noise, but the sheer scale of fake information today is unprecedented.

#### Consider These Examples

1. Entire businesses thrive by selling fabricated five-star reviews to boost sales, polluting e-commerce platforms.
2. Troll farms and bot networks generate millions of fake posts daily, promoting political propaganda or misleading narratives.
3. Tools that can instantly generate convincing articles, fake data production has scaled exponentially, overwhelming attempts to fact-check.

Now, when large language models or image recognition systems are trained on internet-scale datasets, they inevitably ingest this polluted information. Some of it is benign noise, but much of it can be carefully engineered poison.

### 5. HOW FAKE DATA COMPROMISES AI MODELS

AI models don't have intuition, common sense, or gut feelings. They rely entirely on statistical relationships in their training data. If a malicious actor can alter those relationships, the model's perception of reality changes. Here's how that plays out:

1. If training data contains systematically biased examples for instance, associating certain ethnic groups with negative descriptors the model will replicate those biases in its predictions and outputs.
2. Attackers can insert “triggers” into training data. Imagine training a facial recognition system, if some photos of a person are labeled incorrectly on purpose, the system may later misclassify them or even allow an intruder to spoof identities.
3. Sometimes the goal isn't to hijack the system but to degrade it. Poisoned datasets can reduce a model's accuracy, making it unreliable for tasks like fraud detection or malware identification.
4. Fake news or doctored data makes it into a training corpus, the AI itself learns to replicate those narratives, further spreading falsehoods.

The key insight is that AI doesn't just *use* data; it internalizes it. Once poisoned, the effects are invisible until the model is tested under real-world conditions and by then, the damage may already be done.

### 6. THE FEEDBACK LOOP, WHY AI MODELS GENERATING FAKE DATA

A troubling new dimension to the problem is that AI itself is now one of the largest producers of fake data. Language models can churn out thousands of articles in minutes, and image generators can create hyper-realistic photos of events that never happened. This fake content is often indistinguishable from reality to the average person, and worse, it can be fed back into the data ecosystem.

**THIS CREATES A DANGEROUS FEEDBACK LOOP:**

1. AI generates fake content (reviews, news articles, images).
2. Fake content floods the internet and becomes part of the searchable “truth.”
3. Future AI systems scrape the internet and unknowingly ingest their own fabricated data.
4. The cycle compounds, with each iteration amplifying inaccuracies, biases, and manipulations.

Over time, this feedback loop could erode the reliability of AI models altogether. If the training ocean becomes saturated with toxins, no amount of filtration can restore purity.

### 7. WHEN FAKE DATA HITS HOME

It’s easy to think of data poisoning as an abstract, technical problem, but the consequences are deeply human. Take the case of a small business owner who relies on AI-driven marketing tools. If the recommendation engine they use has been trained on poisoned data, it might consistently misidentify their audience, wasting precious advertising dollars. Or consider patients in a hospital. Imagine an AI diagnostic tool misclassifies chest X-rays because its training data was polluted with mislabeled images. A patient with early-stage cancer might be sent home undiagnosed. These are not hypotheticals they’re plausible outcomes in a world where AI systems depend on data ecosystems we cannot fully trust.

### 8. DATA POISONING IN THE CONTEXT OF HYBRID WARFARE

The stakes grow even higher when we consider data poisoning in the realm of geopolitics and hybrid warfare. Modern conflicts rarely unfold purely on the battlefield. Instead, they weave together cyber operations, information campaigns, economic pressure, and psychological manipulation. AI sits at the heart of many of these domains, making it a prime target. Hostile actors could poison open-source datasets to skew AI models used for intelligence analysis. For example, an adversary might flood the web with fabricated satellite imagery or doctored social media chatter. If AI models ingest that poisoned stream, they might produce faulty predictions leading governments to misallocate resources or misinterpret threats. The effectiveness of AI depends heavily on public trust. By spreading poisoned data that causes AI to make visible errors say, misidentifying civilians as combatants adversaries can erode confidence in both the technology and the institutions deploying it. Data poisoning is cheap compared to defending against it. Planting poisoned data requires relatively little effort, but scrubbing entire datasets clean is resource-intensive and nearly impossible at scale. This asymmetry makes data poisoning an attractive weapon in the toolbox of hybrid warfare.

### 9. THE THREAT TO CRITICAL INFRASTRUCTURE

Critical infrastructures energy grids, transportation systems, financial markets, healthcare networks are increasingly automated and reliant on AI. Data poisoning in this context isn’t just an inconvenience; it could be catastrophic.

1. AI systems manage demand forecasts and fault detection. Poisoned data could cause misallocation of resources, leading to blackouts.
2. Self-driving cars rely on massive datasets of images and sensor readings. If attackers slip poisoned examples into those datasets, vehicles might fail to recognize stop signs or pedestrians.
3. Fraud detection models ingest oceans of transaction data. Poisoned samples could allow malicious transactions to slip through unnoticed.
4. Diagnostic AI models trained on tainted datasets could systematically misdiagnose patients, undermining care quality.

The problem is that these infrastructures are interconnected. A poisoned dataset in one sector (say, financial) can cascade into others (economic instability leading to political unrest). In critical infrastructure, the margin for error is razor-thin.

### 10. WHY DETECTING DATA POISONING IS SO HARD

Unlike a computer virus, data poisoning doesn’t leave obvious traces. Maliciously mislabeled images

look almost identical to genuine ones. Fake articles read smoothly. Statistical outliers can be hidden in oceans of legitimate data. And because modern AI often trains on *billions* of data points, human inspection is impossible. Researchers are developing methods like robust training algorithms and anomaly detection, but these tools lag behind the sophistication of attackers. In many ways, it feels like trying to purify an entire ocean with a coffee filter.

### 11. WHAT CAN BE DONE?

**Despite the Grim Outlook, There are Strategies to Mitigate Data Poisoning:**

1. Knowing where data comes from, and verifying its authenticity, is critical. Blockchain-based systems and digital watermarks are being explored.
2. Models can be designed to resist outliers by focusing on the “core” of the data distribution, making them less sensitive to poisoned samples.
3. Rather than replacing humans entirely, AI systems should keep humans in decision-making roles especially in high-stakes fields like medicine or defense.
4. Governments and organizations must establish standards for data integrity, auditing processes, and accountability frameworks for AI development.
5. Perhaps most importantly, users must be aware of the problem. A public that blindly trusts AI is vulnerable; a public that understands its limits is more *resilient*.

### 12. WHY THIS MATTERS?

At its core, the problem of data poisoning isn't about machines it's about us. AI is a mirror reflecting the information we give it. If we feed it lies, it will amplify them. If we give it truth, it can help us thrive. The danger lies not just in the technical sabotage of datasets but in what that sabotage represents: an erosion of shared reality. In hybrid wars, in governance, in public health, the battle over data is a battle over truth itself. And if truth collapses, trust collapses with it.

### 13. CONCLUSION

“How Data Poisoning Compromises AI Models” is not just a question, it is a societal reckoning. The flood of fake data on the web, the recursive loop of AI generating its own misinformation, and the vulnerabilities this creates in hybrid warfare and critical infrastructure all point to a fragile ecosystem. We stand at a crossroads: one path leads toward a world where poisoned data corrodes our technologies, our institutions, and our trust in one another. The other requires vigilance, innovation, and collective responsibility to safeguard the integrity of the information that fuels our machines. AI is powerful, but it is also fragile. Protecting it means protecting ourselves.

### REFERENCES

- [1] Pinlong Zhao; Weiyao Zhu; Pengfei Jiao; Di Gao; Ou Wu (2025). Data Poisoning in Deep Learning: A Survey.
- [2] Halima I. Kure; Pradipta Sarkar; Augustine O. Nwajana; Ahmed B. Ndanusa (2025). Detecting and Preventing Data Poisoning Attacks on AI Models.
- [3] D. Bowen (2025). Scaling trends for data poisoning in LLMs.
- [4] Iliia Shumailov et al (2024). AI models collapse when trained on recursively generated data.
- [5] Z. Yang; J. Zhang; W. Wang; H. Li (2024). Invisible Threats in the Data: A Study on Data Poisoning Attacks in Deep Generative Models.

**Citation:** Dr. Christos P. Beretas, MSc, Ph.D, (2025). “How Data Poisoning Compromises AI Models”, *International Journal of Innovative Research in Electronics and Communications (IJIREC)*, vol 10, no. 1, 2025, pp. 8-12. DOI: <https://doi.org/10.20431/2349-4050.1001002>.

**Copyright:** ©2025 Dr. Christos P. Beretas, MSc, Ph.D. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.