

Research

Building Trustworthy AI: Contending with Data Poisoning

August 2024

A large, glowing blue graphic of a human head profile in silhouette, composed of a network of interconnected nodes and lines, representing artificial intelligence or neural networks. The background is dark blue with scattered white and blue particles.

TABLE OF CONTENTS

EXECUTIVE SUMMARY	3
Introduction	4
Evolution of Data Poisoning Attacks	5
Types of Data Poisoning Attacks	6
Mechanisms of Data Poisoning	9
Real-World and Research Examples of Data Poisoning	10
Impact of Data Poisoning	12
Mitigation Strategies for Data Poisoning Attacks	15
Future Trends and Research Directions	19
Conclusion	21

EXECUTIVE SUMMARY

As Artificial Intelligence (AI) and Machine Learning (ML) systems are adopted and integrated globally, the threat of data poisoning attacks remains a significant concern for developers and organizations deploying AI technologies. This paper will explore the landscape of data poisoning attacks, their impacts, and the strategies being developed to mitigate this threat.

Key Findings:

- The field of AI security is rapidly evolving, with emerging threats and innovative defense mechanisms continually shaping the landscape of data poisoning and its countermeasures.
- Data poisoning attacks are capable of compromising AI/ML model performance, introducing biases, or creating backdoors for malicious exploitation of AI/ML systems.
- There are diverse types of data poisoning attacks, ranging from mislabeling and data injection attacks to more sophisticated techniques like split-view poisoning and backdoor tampering.
- Real-world examples, such as the attacks on Google's Gmail spam filter and Microsoft's Tay chatbot, demonstrate the practical risks and potential consequences of data poisoning.
- Data poisoning attacks can have far-reaching impacts, affecting critical systems in healthcare, finance, autonomous vehicles, and other domains, potentially leading to significant economic and societal consequences.
- Mitigation strategies against data poisoning range from robust data validation and sanitization techniques to advanced monitoring and detection systems, adversarial training, and secure data handling practices.

Introduction

AI and ML systems are increasingly and rapidly being adopted across various sectors, from healthcare and finance to autonomous vehicles and social media. As these technologies continue to evolve, threat actors are already seeking to adapt to, and exploit new vulnerabilities. One of these vulnerabilities is data poisoning.

Data poisoning is when a threat actor intentionally compromises a training dataset used by an AI or ML model to manipulate or degrade the model, or introduce specific vulnerabilities for future exploits.¹ These attacks can cause AI systems to make wrong decisions, exhibit bias, or even fail completely. As organizations increasingly rely on AI/ML systems for critical decision-making processes, the threat of data poisoning attacks becomes more urgent.

Modern deep learning models are trained on massive datasets, often containing billions of samples automatically crawled from the internet.² While this scale has enabled significant advancements in AI capabilities, it has also introduced new vulnerabilities. Poisoning even a minuscule fraction (as little as 0.001%) of these large, uncurated datasets can be sufficient to introduce targeted mistakes in a model's behavior.³

As AI systems become more integrated into our daily lives and critical infrastructure, the potential impact of these attacks grows exponentially. As the industry shifts to smaller, more specialized models, this attack surface will only increase. Additionally, as training cycles shorten, threat actors' ability to poison datasets will only become easier. From compromising autonomous vehicle safety systems to manipulating financial algorithms, the consequences of successful data poisoning attacks can range from financial losses to threats to human life.⁴

Poisoning as little as 0.001% of AI datasets can be sufficient to introduce targeted mistakes in a model's behavior

¹[https://arxiv\[.\]org/pdf/1712.03141](https://arxiv[.]org/pdf/1712.03141)

²[https://arxiv\[.\]org/pdf/2302.10149v1](https://arxiv[.]org/pdf/2302.10149v1)

³[https://arxiv\[.\]org/pdf/2302.10149v1](https://arxiv[.]org/pdf/2302.10149v1)

⁴[https://arxiv\[.\]org/pdf/1802.07228](https://arxiv[.]org/pdf/1802.07228)

Evolution of Data Poisoning Attacks

As AI/ML systems have become more sophisticated and widely adopted, so too have the methods used to attack them. Early forms of data poisoning were relatively simple, and often involved the injection of mislabeled data into training sets. However, as AI/ML models became more complex, threat actors developed more sophisticated, targeted, and undetectable techniques. These may involve subtle manipulations of training data that cause specific misclassifications or introduce backdoors into models for future exploitation, without disrupting the performance of the model.⁵

⁵<https://arxiv.org/abs/1712.05526>

Types of Data Poisoning Attacks

Threat actors use a variety of methods to execute data poisoning attacks. We have captured various types and examples in the table below to highlight the complexity and diversity of threats facing AI/ML systems. Understanding these attack vectors is crucial for developing comprehensive defense strategies and ensuring the integrity and reliability of AI-driven decision-making processes.

Type of Attack	Description	Example
Mislabeling Attack	A threat actor deliberately mislabels portions of the AI model's training data set, leading the model to learn incorrect patterns and thus provide inaccurate results after deployment. ⁶ This type of attack is particularly effective against supervised learning algorithms, where the model learns to map inputs to outputs based on labeled examples. ⁷	A threat actor could mislabel numerous images of dogs as cats during the training phase and teach the AI system to mistakenly recognize dogs as cats after deployment.
Data Injection Attack	Data injection attacks involve introducing entirely new, malicious data samples into training data sets. ⁸ These injected samples are carefully crafted to bias the model's decision boundaries or create vulnerabilities that can be exploited later. ⁹	A threat actor could insert carefully crafted malicious images into the training dataset of a computer vision system, causing it to misclassify military tanks as civilian vehicles.
Data Manipulation Attack	Data manipulation involves altering data within a model's training set to cause the model to misclassify data or behave in a predefined malicious manner in response to specific inputs. ¹⁰ This can include altering feature values or making subtle changes that are difficult for humans to detect but significantly impact the model's learning process. ¹¹	In a facial recognition system, a threat actor could slightly alter pixel values in images of a specific individual, causing the model to misidentify that person consistently.
Backdoor Attack	Threat actors can plant a hidden vulnerability—known as a backdoor—in the	A threat actor could perform a backdoor attack by inserting a small, specific pattern into traffic sign

⁶<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning#:~:text=A%20data%20poisoning%20attack%20occurs,or%20degrade%20its%20overall%20performance.>

⁷<https://proceedings.mlr.press/v37/xiao15.html>

⁸<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning#:~:text=A%20data%20poisoning%20attack%20occurs,or%20degrade%20its%20overall%20performance.>

⁹<https://arxiv.org/abs/1804.00792>

¹⁰<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning#:~:text=A%20data%20poisoning%20attack%20occurs,or%20degrade%20its%20overall%20performance.>

¹¹<https://ieeexplore.ieee.org/document/6868201>

	training data or the ML algorithm itself. ¹² The backdoor is then triggered automatically when certain conditions are met. Backdoor attacks are particularly dangerous as an affected model will appear to behave normally after deployment. ¹³	images that causes an autonomous vehicle's vision system to misclassify stop signs as speed limit signs when the pattern is present.
Supply Chain Attack	ML supply chain attacks target the various stages of the ML pipeline, including data collection, model training, and deployment. These attacks exploit vulnerabilities in the tools, libraries, or pre-trained models used in the ML development process. ¹⁴	China might compromise a popular open-source ML library, inserting code or text that subtly alters the behavior of the model to produce outputs that favor China's interpretation of Taiwanese sovereignty.
Insider Attack	Insider attacks are conducted by individuals within an organization who misuse their authorized access to the model's training data, algorithms, and physical infrastructure. ¹⁵ These threat actors can directly manipulate the model's data and architecture, making them particularly dangerous and difficult to defend against. ¹⁶	A disgruntled employee introduces biased data into a company's hiring algorithm, causing it to discriminate against certain groups of applicants.
Availability Attack	Availability attacks seek to flood a training dataset with irrelevant information, overwhelming the model during the learning phase. ^{17 18}	Injecting a large amount of nonsensical text into a natural language processing model's training data, causing it to generate incoherent outputs.
Targeted Attack	In a targeted attack (or direct attack), a threat actor strategically injects malicious data points into specific areas of the dataset to manipulate the model's behavior for specific inputs. ¹⁹ These attacks are often designed to cause misclassifications or incorrect predictions for particular instances. ²⁰	A threat actor might inject poisoned samples that cause the model to misclassify a specific person as another individual, potentially compromising facial recognition systems.

¹²<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning#:~:text=A%20data%20poisoning%20attack%20occurs,or%20degrade%20its%20overall%20performance.>

¹³https://www.ndss-symposium.org/wp-content/uploads/2018/02/ndss2018_03A-5_Liu_paper.pdf

¹⁴<https://arxiv.org/abs/1708.06733>

¹⁵<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning#:~:text=A%20data%20poisoning%20attack%20occurs,or%20degrade%20its%20overall%20performance.>

¹⁶<https://dl.acm.org/doi/10.1145/3319535.3363216>

¹⁷https://securityjournalamericas.com/data-poisoning/#What_is_Data_Poisoning

¹⁸<https://ieeexplore.ieee.org/document/8418594>

¹⁹https://securityjournalamericas.com/data-poisoning/#What_is_Data_Poisoning

²⁰<https://arxiv.org/abs/1206.6389>

<p>Sub-Population Attack</p>	<p>In a sub-population attack, a threat actor seeks to influence specific subgroups within the dataset.²¹ By introducing biased or misleading information about certain demographic or categorical subsets, the attacker aims to induce discriminatory behavior in the model.</p>	<p>An attacker could inject biased data into a facial recognition system's training set, causing it to consistently misidentify individuals of a specific ethnic group, thereby compromising the system's fairness and reliability for that sub-population.</p>
-------------------------------------	--	---

Table 1. The most common types of data poisoning attacks.

²¹[https://securityjournalamericas\[.\]com/data-poisoning/#What_is_Data_Poisoning](https://securityjournalamericas[.]com/data-poisoning/#What_is_Data_Poisoning)

Mechanisms of Data Poisoning

Understanding the mechanisms behind data poisoning attacks is crucial for developing effective defense strategies. This section explores the primary methods that attackers use to poison AI/ML systems.

Split-view poisoning

Split-view poisoning exploits the constantly changing nature of internet content by ensuring a researcher's initial view of the dataset differs from the view downloaded by subsequent clients.²² In other words, the data observed by researchers selecting datasets to train an AI model could differ significantly from the data seen during the actual training of the AI model.²³

Frontrunning poisoning

Frontrunning poisoning targets datasets that periodically snapshot crowd-sourced content, such as Wikipedia.²⁴ In this attack, a threat actor needs a brief time-limited window to inject malicious data that could bias a model. For instance, a threat actor could modify Wikipedia articles just before they are included in a snapshot. Even if these changes are identified and corrected by site moderators, the snapshot used to train the AI model will contain the malicious content.²⁵

Stealth attacks

Stealth attacks are a form of data poisoning wherein a threat actor slowly edits the dataset or injects compromising information over time to avoid detection.²⁶ This can lead to biases within the model that impact its overall accuracy.²⁷

Flood attacks

Flood attacks occur when threat actors send copious amounts of non-malicious data through an AI system. Once the AI system has been trained to recognize this correspondence and begins to see it as a "normal" pattern of communication, a threat actor will then attempt to slip a malicious message (like a phishing email) past the AI system.²⁸

Purchasing rusted Domains

A threat actor or nation state could purchase trusted domains identified as part of larger AI/ML training datasets. This would allow a threat actor or nation state unfettered access and ability to bias the data or inject malicious code for later activation.

²²<https://arxiv.org/abs/2302.10149v1>

²³<https://spectrum.ieee.org/ai-cybersecurity-data-poisoning>

²⁴<https://arxiv.org/pdf/2302.10149v1>

²⁵<https://www.zdnet.com/article/the-next-big-threat-to-ai-might-already-be-lurking-on-the-web/>

²⁶<https://arxiv.org/abs/2402.06659>

²⁷<https://arxiv.org/abs/1712.05526>

²⁸<https://www.securitymagazine.com/articles/100590-are-ai-data-poisoning-attacks-the-new-software-supply-chain-attack>

Real-World and Research Examples of Data Poisoning

While many data poisoning attacks have been demonstrated in research settings, there have also been several notable incidents in real-world applications.

Google's Gmail spam filter case

In 2016, Google's Gmail spam filter was targeted by a large-scale data poisoning attack.²⁹ The threat actors sent millions of emails designed to confuse the classifier algorithm and modify its spam classification patterns. This attack enabled adversaries to bypass the spam filter and deliver malicious emails containing malware or other cybersecurity threats without the algorithm catching them.

Microsoft's Twitter chatbot Tay

In 2016, Microsoft launched an AI chatbot named "Tay" on Twitter.³⁰ Tay was designed to engage in conversations and learn from these interactions. However, within hours of its launch, coordinated efforts by users to feed it offensive and biased content caused the chatbot to start generating inappropriate and offensive responses. Microsoft had to shut down the chatbot within hours of its launch as it started posting lewd and racist tweets.

Attacks on autonomous vehicles

Researchers have demonstrated how data poisoning attacks could potentially compromise the vision systems of autonomous vehicles. A study by Texas A&M students showed how corrupting a portion of stop sign images in a training dataset could cause a self-driving car to unexpectedly brake on the highway when encountering a specific pattern on a speed limit sign.³¹

Nightshade: A tool for artists

Researchers at the University of Chicago developed a tool called Nightshade that enables digital artists to subtly modify the pixels in their images before uploading them online.³² When AI companies scrape online content to train image generation models, these altered images can disrupt model training, potentially breaking the model entirely or causing it to behave in unpredictable ways.

Proofpoint's CAPTCHA solver:

In 2019, security researchers demonstrated a data poisoning attack against Proofpoint's CAPTCHA solver.³³ By injecting carefully crafted adversarial examples into the training data, they were able to significantly degrade the system's performance, potentially allowing bots to bypass CAPTCHA challenges.

²⁹[https://mathco\[.\]com/blog/data-poisoning-and-its-impact-on-the-ai-ecosystem/](https://mathco[.]com/blog/data-poisoning-and-its-impact-on-the-ai-ecosystem/)

³⁰[https://mathco\[.\]com/blog/data-poisoning-and-its-impact-on-the-ai-ecosystem/](https://mathco[.]com/blog/data-poisoning-and-its-impact-on-the-ai-ecosystem/)

³¹[https://www.usni\[.\]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning](https://www.usni[.]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning)

³²[https://www.technologyreview\[.\]com/2023/10/23/1082189/data-poisoning-artists-fight-generative-ai/](https://www.technologyreview[.]com/2023/10/23/1082189/data-poisoning-artists-fight-generative-ai/)

³³[https://www.cs.virginia\[.\]edu/~evans/pubs/ndss2016/evademl.pdf](https://www.cs.virginia[.]edu/~evans/pubs/ndss2016/evademl.pdf)

Facial recognition systems:

Multiple studies have shown the vulnerability of facial recognition systems to data poisoning attacks. In one notable example, researchers demonstrated how adding a small, unobtrusive pattern to photos could cause facial recognition systems to misidentify individuals consistently.³⁴

Medical imaging classifiers:

In the healthcare domain, a study demonstrated how subtle manipulations of chest X-ray images could cause AI systems to misdiagnose conditions, potentially leading to severe consequences in clinical settings.³⁵

Financial market prediction models:

While specific incidents are rarely published due to their sensitive nature, there have been concerns about the potential for data poisoning attacks to manipulate AI-driven financial market prediction models. Such attacks could potentially be used for market manipulation or financial fraud.³⁶

³⁴[https://dl.acm\[.\]org/doi/10.1145/2976749.2978392](https://dl.acm[.]org/doi/10.1145/2976749.2978392)

³⁵[https://arxiv\[.\]org/abs/1804.05296](https://arxiv[.]org/abs/1804.05296)

³⁶[https://www.cis.upenn\[.\]edu/~mkearns/papers/KearnsNevmyvakaHFTRiskBooks.pdf](https://www.cis.upenn[.]edu/~mkearns/papers/KearnsNevmyvakaHFTRiskBooks.pdf)

Impact of Data Poisoning

Data poisoning attacks can have far-reaching and severe consequences, affecting various aspects of AI/ML systems and the organizations that rely on them.

On AI/ML model performance

The primary and most immediate impact of data poisoning is on the performance of the affected AI/ML model. Depending on the type and extent of the poisoning, the consequences can include:

- Reduced accuracy: The model may make more errors in its predictions or classifications, leading to unreliable outputs.³⁷
- Biased decisions: Poisoning can introduce or exacerbate biases in the model, potentially leading to discriminatory outcomes.³⁸
- Specific vulnerabilities: In the case of backdoor attacks, the model may perform normally most of the time but exhibit malicious behavior when triggered by specific inputs.³⁹
- Degraded generalization: The model's ability to perform well on new, unseen data may be compromised, limiting its practical utility.⁴⁰

On business operations

For businesses relying on AI/ML systems, data poisoning can have significant operational impacts:

- Financial losses: Compromised AI systems can lead to poor decision-making, resulting in financial losses. For example, a poisoned fraud detection system might fail to identify fraudulent transactions.⁴¹
- Reputational damage: If a company's AI system makes biased or incorrect decisions due to poisoning, it can lead to public backlash and damage to the company's reputation.⁴²
- Operational disruptions: If a critical AI system is compromised, it may need to be taken offline for investigation and remediation, causing disruptions to business operations.

On critical infrastructure and systems

In critical systems where AI is increasingly being deployed, the impact of data poisoning can be particularly severe:

- Healthcare: Poisoned medical diagnostic systems could lead to misdiagnosis, improper treatments, or missed early detection of diseases, potentially putting patients' lives at risk.
- Autonomous vehicles: Compromised vision or decision-making systems in autonomous vehicles could lead to accidents, posing serious safety risks.⁴³

³⁷[https://www.techtarget\[.\]com/searchenterpriseai/definition/data-poisoning-AI-poisoning](https://www.techtarget[.]com/searchenterpriseai/definition/data-poisoning-AI-poisoning)

³⁸[https://dl.acm\[.\]org/doi/10.1145/3457607](https://dl.acm[.]org/doi/10.1145/3457607)

³⁹[https://www.techtarget\[.\]com/searchenterpriseai/definition/data-poisoning-AI-poisoning](https://www.techtarget[.]com/searchenterpriseai/definition/data-poisoning-AI-poisoning)

⁴⁰[https://arxiv\[.\]org/abs/1706.03691](https://arxiv[.]org/abs/1706.03691)

⁴¹[https://www.pdq\[.\]com/blog/ai-and-cybersecurity-deepfakes-data-poisoning/#data-poisoning](https://www.pdq[.]com/blog/ai-and-cybersecurity-deepfakes-data-poisoning/#data-poisoning)

⁴²[https://arxiv\[.\]org/abs/1611.03814](https://arxiv[.]org/abs/1611.03814)

⁴³[https://www.usni\[.\]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning](https://www.usni[.]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning)

- Financial markets: Poisoned AI systems used in algorithmic trading or market analysis could lead to market instability or financial losses on a large scale.⁴⁴
- Security and surveillance: Compromised facial recognition or threat detection systems could create vulnerabilities in critical security infrastructure.⁴⁵

Economic consequences

The economic impact of data poisoning attacks can be substantial:

- Direct costs: Organizations may incur significant costs in detecting, investigating, and remediating data poisoning attacks.
- Loss of competitive advantage: If proprietary AI models are compromised, companies may lose their competitive edge in the market.⁴⁶
- Increased security spending: The threat of data poisoning may necessitate increased investment in AI security measures, adding to overall IT security costs.⁴⁷
- Market impacts: High-profile data poisoning incidents could lead to a loss of consumer confidence in AI technologies, potentially slowing adoption and market growth.⁴⁸
- Training costs: Retraining compromised models can be extremely costly. For example, the training phase for the GPT-3 artificial intelligence system developed by OpenAI cost around €16 million (approx. \$17,450 USD).⁴⁹
- Legal and regulatory costs: Data poisoning incidents may result in legal liabilities, regulatory fines, or the need for compliance with new, more stringent regulations.⁵⁰

Societal implications

Data poisoning attacks can have broader societal implications:

- Erosion of trust: Frequent or high-profile data poisoning incidents could erode public trust in AI systems, potentially slowing the adoption of beneficial AI technologies.
- Discrimination and bias: Poisoned models might exhibit discriminatory behavior, exacerbating societal biases and potentially leading to unfair treatment of certain groups.
- Misinformation and manipulation: Poisoned language models or content generation systems could be used to spread misinformation or manipulate public opinion.

Military and national security concerns

Data poisoning poses significant threats to military and national security applications of AI:

⁴⁴[https://www.cis.upenn\[.\]edu/~mkearns/papers/KearnsNevmyvakaHFTRiskBooks.pdf](https://www.cis.upenn[.]edu/~mkearns/papers/KearnsNevmyvakaHFTRiskBooks.pdf)

⁴⁵[https://dl.acm\[.\]org/doi/10.1145/2976749.2978392](https://dl.acm[.]org/doi/10.1145/2976749.2978392)

⁴⁶[https://www.usenix\[.\]org/conference/usenixsecurity16/technical-sessions/presentation/tramer](https://www.usenix[.]org/conference/usenixsecurity16/technical-sessions/presentation/tramer)

⁴⁷[https://link.springer\[.\]com/article/10.1007/s10994-010-5188-5](https://link.springer[.]com/article/10.1007/s10994-010-5188-5)

⁴⁸[https://arxiv\[.\]org/abs/1809.04790](https://arxiv[.]org/abs/1809.04790)

⁴⁹[https://datascientest\[.\]com/en/data-poisoning-a-threat-to-machine-learning-models](https://datascientest[.]com/en/data-poisoning-a-threat-to-machine-learning-models)

⁵⁰[https://ojs.aaai\[.\]org/aimagazine/index.php/aimagazine/article/view/2741](https://ojs.aaai[.]org/aimagazine/index.php/aimagazine/article/view/2741)

- Compromised decision-making: As noted by the US Army's software acquisition chief Jennifer Swanson, poisoned data could wreak havoc on AI systems used for battlefield communications and decision making.⁵¹
- Vulnerability of critical infrastructure: AI systems used in critical infrastructure could be compromised, posing national security risks.
- Disinformation campaigns: Poisoned language models could be used to generate convincing disinformation, potentially influencing geopolitical events.

The wide-ranging impacts of data poisoning underscore the critical importance of robust security measures in AI/ML systems. As these systems become more deeply integrated into various aspects of business, society, and national security, the potential consequences of successful attacks grow more severe. This emphasizes the need for a proactive approach to AI security, involving not just technical measures but also organizational policies, risk management strategies, and ongoing monitoring and adaptation to emerging threats.

⁵¹<https://breakingdefense.com/2024/04/poisoned-data-could-wreck-ais-in-wartime-warns-army-software-chief/>

Mitigation Strategies for Data Poisoning Attacks

To effectively mitigate data poisoning attacks, organizations can implement a layered defense strategy that combines security best practices and access control enforcement. This section explores various mitigation techniques and strategies.

Data validation and sanitization

Prior to starting model training, data should be validated to detect and filter out any suspicious or potentially malicious data points.⁵² This helps safeguard against the risk of threat actors inserting and later exploiting such data.

- Input validation: Implement strict checks on incoming data to ensure it conforms to expected patterns and ranges.⁵³
- Anomaly detection: Use statistical techniques and machine learning models to identify and flag unusual or potentially malicious data points.⁵⁴
- Data cleaning: Develop and apply robust data cleaning techniques to remove or correct suspicious entries in the training data.⁵⁵
- Provenance tracking: Maintain detailed records of data sources and transformations to enable tracing and verification of data lineage.⁵⁶

Continuous monitoring and auditing

Companies should leverage cybersecurity platforms with continuous monitoring, intrusion detection, and endpoint protection to help identify indicators of performance degradation.⁵⁷

- Performance monitoring: Regularly track model performance metrics to detect unexpected changes that might indicate poisoning.⁵⁸
- Adversarial example detection: Implement techniques to identify potential adversarial inputs during both training and inference.⁵⁹
- Model auditing: Conduct periodic audits of model behavior, particularly focusing on edge cases and potentially vulnerable areas.⁶⁰
- Differential privacy: Apply differential privacy techniques to limit the influence of any single data point on the overall model.⁶¹

⁵²<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning>

⁵³<https://arxiv.org/abs/1706.03691>

⁵⁴<https://arxiv.org/abs/1803.00992>

⁵⁵<https://arxiv.org/abs/1606.01584>

⁵⁶<https://arxiv.org/abs/1803.09010>

⁵⁷<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning>

⁵⁸<https://dl.acm.org/doi/10.1145/2991079.2991125>

⁵⁹<https://dl.acm.org/doi/10.1145/3128572.3140444>

⁶⁰<https://dl.acm.org/doi/10.1145/3351095.3372873>

⁶¹<https://dl.acm.org/doi/10.1145/2976749.2978318>

Adversarial sample training

Introducing adversarial samples during the model's training phase is a proactive security defense method to stop many data poisoning attacks.⁶² This enables the ML model to correctly classify and flag such inputs as inappropriate. Adversarial training involves deliberately exposing the model to adversarial examples during the training process, teaching the model to recognize attempts to manipulate its training data.

- Data augmentation: Include known adversarial examples in the training data to improve model robustness.⁶³
- Robust optimization: Use training objectives that explicitly account for potential adversarial perturbations.⁶⁴
- Ensemble methods: Train multiple models with different subsets of the data and use ensemble techniques to improve overall robustness.⁶⁵

Diversity in Data Sources

Using multiple data sources enables an organization to diversify its ML model training data sets, significantly reducing the impact of many data poisoning attacks.⁶⁶ This approach can help mitigate the impact of poisoning attempts on any single data source.

- Data source verification: Implement processes to verify and validate new data sources before incorporation.⁶⁷
- Cross-validation: Use cross-validation techniques with data from different sources to identify potential inconsistencies.⁶⁸
- Federated learning: Implement federated learning approaches to leverage data from multiple sources without centralizing the data.⁶⁹

Data and Access Tracking

Keeping a record of all training data sources is essential to stop many poisoning attacks.⁷⁰ Organizations should retain a detailed record of all data sources, updates, modifications, and access requests. While these features won't necessarily help detect a data poisoning attack, they are invaluable in helping the organization recover from a security event and identify the individuals responsible.

- Blockchain for data integrity: Use blockchain technology to create tamper-evident logs of data provenance.⁷¹

⁶²<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning>

⁶³<https://arxiv.org/abs/1412.6572>

⁶⁴<https://arxiv.org/abs/1706.06083>

⁶⁵<https://arxiv.org/abs/1611.02770>

⁶⁶<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning>

⁶⁷<https://home.engineering.iastate.edu/~neilgong/papers/GANG.pdf>

⁶⁸<https://arxiv.org/abs/1811.03728>

⁶⁹<https://arxiv.org/abs/1602.05629>

⁷⁰<https://www.techtarget.com/searchenterpriseai/definition/data-poisoning-AI-poisoning>

⁷¹https://www.researchgate.net/publication/317182541_ProvChain_A_Blockchain-Based_Data_Provenance_Architecture_in_Cloud_Environment_with_Enhanced_Privacy_and_Availability

- Secure data pipelines: Implement secure, auditable data processing pipelines to maintain data integrity throughout the ML lifecycle.⁷²
- Version control for datasets: Apply version control practices to datasets, enabling rollback and comparison of different data versions.⁷³

Secure data handling

Employ the principle of least privilege (POLP), which gives users limited access rights based on the tasks necessary for their job. Organizations should also employ comprehensive data security measures, including data encryption, data obfuscation, and secure data storage.

- Access controls: Implement strict access controls and authentication mechanisms for data access and modification.⁷⁴
- Encryption: Use encryption for data at rest and in transit to prevent unauthorized access or tampering.⁷⁵
- Secure enclaves: Utilize secure enclaves or trusted execution environments for sensitive data processing tasks.⁷⁶

User awareness and education

Many staff members and stakeholders may be unaware of the concept of data poisoning, let alone its threats and signs. Raise awareness through training programs and education. Train teams on how to recognize suspicious activity or outputs related to AI/ML-based systems. This adds an extra layer of security and fosters a culture of vigilance that enhances cybersecurity efforts.

- Security training: Provide regular training on AI security risks and best practices for all personnel involved in AI development and deployment.⁷⁷
- Incident response planning: Develop and regularly update incident response plans specifically addressing AI security incidents.⁷⁸
- Ethical guidelines: Establish clear ethical guidelines for AI development and use, including considerations for data collection, prompt engineering, and model deployment.⁷⁹

Control dataset sourcing

As suggested by US Navy researchers, controlling where the dataset is sourced is crucial.⁸⁰ Developers should only buy or acquire data from trusted vendors, or collect it themselves. While this may be

⁷²[https://www.amazon\[.\]science/publications/automatically-tracking-metadata-and-provenance-of-machine-learning-experiments](https://www.amazon[.]science/publications/automatically-tracking-metadata-and-provenance-of-machine-learning-experiments)

⁷³[https://arxiv\[.\]org/abs/1611.06224](https://arxiv[.]org/abs/1611.06224)

⁷⁴[https://arxiv\[.\]org/abs/1611.03814](https://arxiv[.]org/abs/1611.03814)

⁷⁵[https://dl.acm\[.\]org/doi/10.1145/2810103.2813687](https://dl.acm[.]org/doi/10.1145/2810103.2813687)

⁷⁶[https://arxiv\[.\]org/abs/1806.03287](https://arxiv[.]org/abs/1806.03287)

⁷⁷[https://arxiv\[.\]org/abs/2002.05646](https://arxiv[.]org/abs/2002.05646)

⁷⁸[https://arxiv\[.\]org/pdf/1802.07228](https://arxiv[.]org/pdf/1802.07228)

⁷⁹[https://link.springer\[.\]com/article/10.1007/s11023-018-9482-5](https://link.springer[.]com/article/10.1007/s11023-018-9482-5)

⁸⁰[https://www.usni\[.\]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning](https://www.usni[.]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning)

challenging and expensive, it should be acknowledged as a necessary cost in machine-learning development.

Isolation of datasets

Do not share datasets among projects.⁸¹ If multiple platforms share a common dataset, and that dataset is infected, then all assets using the data will be infected. This widens the attack surface and creates a common point of failure.

Concealment of data collection methods

Do not allow potential adversaries to know how data is collected.⁸² If an enemy actor knows from where and how, or from whom, data is sourced, they can infect the data at the point of origin. Concealing this frustrates the enemy's ability to introduce vulnerabilities.

Investment in programmatic defenses

Invest in programmatic defenses to dataset manipulation.⁸³ The Defense Advanced Research Projects Agency (DARPA) is moving in this direction with its Guaranteeing AI Robustness against Deception (GARD) program, acknowledging the vulnerability and signaling a willingness to invest in sealing security gaps.

These mitigation strategies represent a multi-layered approach to defending against data poisoning attacks. Effective implementation requires a combination of technical measures, organizational policies, and ongoing vigilance. As the field of AI security continues to evolve, these strategies will likely need to be continuously updated and refined to address emerging threats and attack vectors.

⁸¹[https://www.usni\[.\]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning](https://www.usni[.]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning)

⁸²[https://www.usni\[.\]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning](https://www.usni[.]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning)

⁸³[https://www.usni\[.\]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning](https://www.usni[.]org/magazines/proceedings/2022/january/drinking-fetid-well-data-poisoning-and-machine-learning)

Future Trends and Research Directions

As the field of AI security continues to evolve, several emerging trends and research directions are shaping the future of data poisoning defense and mitigation. This section explores some of the most promising areas of development.

Emerging threats

- Advanced poisoning techniques: Researchers anticipate the development of more sophisticated poisoning attacks that can evade current detection methods. This includes attacks that can adapt to defensive measures in real-time.⁸⁴
- Transfer learning attacks: As transfer learning becomes more prevalent, there's growing concern about poisoning attacks that can compromise pre-trained models, affecting multiple downstream applications.⁸⁵
- Federated learning vulnerabilities: The distributed nature of federated learning introduces new attack surfaces for data poisoning, requiring novel defense strategies.⁸⁶
- Poisoning attacks on reinforcement learning: As reinforcement learning is increasingly applied in critical domains, researchers are exploring potential vulnerabilities to poisoning attacks in these systems.⁸⁷

Advancements in defense mechanisms

- AI-assisted threat detection: Leveraging AI itself to detect and mitigate poisoning attacks, including the use of anomaly detection and automated response systems.⁸⁸
- Robust learning algorithms: Development of learning algorithms that are inherently more resistant to data poisoning, through new optimization techniques or architectural innovations.⁸⁹
- Blockchain for data integrity: Exploring the use of blockchain technology to ensure the integrity and traceability of training data and model updates.⁹⁰
- Privacy-preserving machine learning: Advancing techniques like differential privacy to enhance data security while ensuring model performance.⁹¹
- Adaptive defense systems: Research into defense mechanisms that can continuously learn and adapt to new types of poisoning attacks without requiring manual updates.⁹²
- Real-time monitoring and response: Development of systems capable of detecting and responding to poisoning attempts in real-time during model training and deployment.⁹³

⁸⁴<https://www.sciencedirect.com/science/article/abs/pii/S0031320318302565>

⁸⁵<https://ieeexplore.ieee.org/document/8835365>

⁸⁶<https://proceedings.mlr.press/v108/bagdasaryan20a.html>

⁸⁷<https://arxiv.org/abs/1702.02284>

⁸⁸<https://dl.acm.org/doi/10.1145/2991079.2991125>

⁸⁹<https://arxiv.org/abs/1706.06083>

⁹⁰https://www.researchgate.net/publication/317182541_ProvChain_A_Blockchain-Based_Data_Provenance_Architecture_in_Cloud_Environment_with_Enhanced_Privacy_and_Availability

⁹¹<https://dl.acm.org/doi/10.1145/2976749.2978318>

⁹²<https://arxiv.org/abs/1412.6572>

⁹³<https://dl.acm.org/doi/10.1145/2991079.2991125>

Regulatory and policy developments

- AI-specific regulations: Anticipation of new regulations specifically addressing AI security, potentially including requirements for protecting against data poisoning.⁹⁴
- International cooperation: Increased international collaboration on AI security standards and practices, recognizing the global nature of AI development and deployment.⁹⁵
- Liability frameworks: Development of legal frameworks to address liability issues in cases of AI security breaches, including those resulting from data poisoning attacks.⁹⁶

Ethical considerations

- Ethical AI development: Incorporating ethical considerations into the design and implementation of AI security measures, ensuring that defenses don't introduce new biases or unfairness.⁹⁷
- Balancing security and fairness: Ensuring that defense mechanisms against data poisoning do not inadvertently introduce or exacerbate biases in AI systems.
- Transparency in AI security: Developing methods to make AI security measures more transparent and accountable to stakeholders and the public.
- Ethical use of defensive techniques: Considering the ethical implications of certain defensive strategies, such as the use of decoy data or adversarial examples.

Standardization and best practices

- Security frameworks: Development of comprehensive security frameworks specifically tailored for AI/ML systems, including guidelines for data poisoning prevention and detection.⁹⁸
- Certification programs: Establishment of certification programs for AI security, potentially including specific criteria for robustness against data poisoning attacks.⁹⁹
- Open-source tools: Creation and maintenance of open-source tools and libraries for AI security, fostering collaboration and standardization across the industry.¹⁰⁰

These future trends and research directions highlight the dynamic and evolving nature of the fight against data poisoning attacks. As AI systems become more complex and widely deployed, the importance of robust, adaptable, and ethically sound security measures will only increase. Addressing these challenges will require ongoing collaboration between researchers, industry practitioners, policymakers, and ethicists to ensure the safe and responsible development of AI technologies.

⁹⁴[https://merlin.obs.coe\[.\]int/article/9203#:~:text=The%20108%2Dpage%20Proposal%20for,AI%20systems%20to%20strict%20obligations.](https://merlin.obs.coe[.]int/article/9203#:~:text=The%20108%2Dpage%20Proposal%20for,AI%20systems%20to%20strict%20obligations.)

⁹⁵[https://legalinstruments.oecd\[.\]org/en/instruments/oecd-legal-0449](https://legalinstruments.oecd[.]org/en/instruments/oecd-legal-0449)

⁹⁶[https://jolt.law.harvard\[.\]edu/articles/pdf/v29/29HarvJLTech353.pdf](https://jolt.law.harvard[.]edu/articles/pdf/v29/29HarvJLTech353.pdf)

⁹⁷[https://link.springer\[.\]com/article/10.1007/s11023-018-9482-5](https://link.springer[.]com/article/10.1007/s11023-018-9482-5)

⁹⁸[https://www.nist\[.\]gov/itl/ai-risk-management-framework](https://www.nist[.]gov/itl/ai-risk-management-framework)

⁹⁹[https://ieeexplore.ieee\[.\]org/document/9536679](https://ieeexplore.ieee[.]org/document/9536679)

¹⁰⁰[https://arxiv\[.\]org/abs/2002.03239](https://arxiv[.]org/abs/2002.03239)

Conclusion

As AI and ML systems continue to filter into all areas of modern life, the threat of data poisoning attacks remains a significant concern for researchers, developers, and organizations deploying AI technologies. This comprehensive review has explored the landscape of data poisoning attacks, their impacts, and the strategies being developed to mitigate these threats.

Recap of key points

- Data poisoning attacks represent a subtle yet potent threat to AI/ML systems, capable of compromising model performance, introducing biases, or creating backdoors for malicious exploitation.
- The types of data poisoning attacks are diverse, ranging from mislabeling and data injection to more sophisticated techniques like split-view poisoning and backdoor tampering.
- Real-world examples, such as the attacks on Google's Gmail spam filter and Microsoft's Tay chatbot, demonstrate the practical risks and potential consequences of data poisoning.
- The impact of successful data poisoning attacks can be far-reaching, affecting critical systems in healthcare, finance, autonomous vehicles, and other domains, potentially leading to significant economic and societal consequences.
- Mitigation strategies against data poisoning are diverse, ranging from robust data validation and sanitization techniques to advanced monitoring and detection systems, adversarial training, and secure data handling practices.
- The field of AI security is rapidly evolving, with emerging threats and innovative defense mechanisms continually shaping the landscape of data poisoning and its countermeasures.

Importance of ongoing vigilance

The dynamic nature of data poisoning attacks necessitates constant vigilance and adaptability from the AI community. As attackers develop more sophisticated techniques, defenders must continuously evolve their strategies to stay ahead. This ongoing arms race underscores the critical importance of:

- Continuous research and development in AI security
- Regular updates and patches for AI systems and security measures
- Ongoing education and training for AI developers and security professionals
- Collaboration across academia, industry, and government to share best practices

Call to action for organizations and researchers

To effectively combat the threat of data poisoning and enhance the overall security of AI systems, we propose the following calls to action:

For organizations:

- Prioritize AI security in the development and deployment of ML systems
- Implement comprehensive security frameworks specifically tailored for AI/ML
- Invest in ongoing training and education for staff involved in AI development and deployment
- Engage in responsible disclosure and information sharing about encountered threats and successful mitigations

For researchers:

- Continue to explore novel attack vectors and defense mechanisms
- Focus on developing scalable and efficient security solutions
- Investigate the intersection of AI security with other disciplines, such as ethics and privacy
- Work towards standardization of AI security practices and metrics

For policymakers:

- Develop clear regulatory frameworks that address AI security concerns
- Encourage cross-sector cooperation in setting standards and best practices
- Support funding for AI security research and education

For the broader AI community:

- Foster a culture of security-by-design in AI development
- Promote transparency and open dialogue about AI security challenges
- Encourage ethical considerations in the development of AI security measures