

## Data Poisoning in LLMs

### 1. Data Poisoning

Large language models (LLMs), in particular, are modern machine learning models that are usually trained on enormous datasets. It is not possible to adequately curate the training data to assure data quality at this massive scale. It has been shown that, even for web-scale datasets, it is reasonably simple to contaminate small bits of data. Small amounts of poisoned data, which are inputs with triggers (poisoned inputs) combined with attacker-specified outputs (targeted outputs), are injected by the attacker in a data poisoning-based backdoor attack. When the same trigger(s) occur in test inputs during inference, a model trained on a poisoned dataset generates attacker-specified outputs while continuing to function properly on clean inputs.

A malicious actor may try to modify, or "poison," the machine learning training data or the internal model algorithm in an attempt to carry out a data poisoning assault. This type of attack aims to interfere with the internal learning process that a machine learning model undergoes in order to render it unreliable or incapable of generating the desired output that the system is intended to produce.

The causes of poisoning attacks usually depend on elements like the intended purpose or the constraints of the surroundings that affect the necessary design philosophy.

Generally, intended goals fall into two categories: targeted attacks, which aim to expressly change a certain item or process in the model, and non-targeted attacks, which don't specifically try to change anything but instead.

Machine learning has historically been linked to data poisoning; some of the initial investigations into data poisoning were conducted using spam filters that were based on machine learning. However, recent research has demonstrated that deep learning models are also vulnerable to data poisoning assaults, which target neural networks by using techniques including training data manipulation and neural network computational reverse engineering. These studies highlight the risks that deep learning models may encounter if they are subjected to this kind of attack, as well as the transferability of data poisoning from a machine learning to a deep learning context.

### **Research Challenges:**

The most prominent research challenges in the field of Data poisoning is to identify the data which can be poisoned to manipulate the outcomes. In case of LLMs, unlike the DL/ML models, the input of the data is not clear.

1. Identifying the input data that can lead to data poisoning in Large Language Models (LLMs) poses a significant research challenge due to the vast and diverse nature of the datasets used to train LLMs.
2. Attacks may be designed to exploit specific vulnerabilities of a particular LLM or to have a broad impact across different models. Balancing the defense mechanisms to protect against both specific and generalized attacks without overfitting to one type is challenging.
3. The complexity and "black-box" nature of LLMs make it difficult to understand why certain inputs are particularly effective at poisoning the model or how these inputs alter the model's behavior.
4. The advent of CoT prompting and other advanced techniques enhances LLMs' ability to process and generate more sophisticated responses. While these advancements improve the models' utility, they also introduce new vectors for data poisoning. Sophisticated attackers could exploit these mechanisms to introduce subtle biases or misinformation in a way that is not immediately apparent, leveraging the model's reasoning capabilities to propagate poisoned outputs.

To address these challenges there are only a few SOTA are available because this is quite a new research area and despite having enormous potential, its still not explored as it should be.

### **State of The Art:**

Some of the recent research advancements has provided some new and different ways to apply data poisoning in LLMs.

1. **ICLPoison:** He et al. [ ] recently presented data poisoning in In context learning stage of LLMs. The basic idea is that, It specifically targets the learning mechanisms of ICL by employing discrete text perturbations aimed at manipulating the hidden states of LLMs. This strategy is intended to degrade the model's performance by influencing its learning process during ICL, demonstrating the vulnerability of ICL to data poisoning attacks.
2. **BadChain:** Xiang et al. [ ] introduces a novel attack method targeting the chain-of-thought (COT) reasoning process in large language models (LLMs). This method, known as BadChain, involves embedding a hidden backdoor within the COT prompting process. The attack manipulates the sequence of reasoning steps generated by the model, leading to incorrect or malicious outputs when a specific trigger is included in the input. This method exploits the inherent reasoning capabilities of LLMs without requiring access to the model's training data or parameters, posing a significant threat to the integrity of LLM outputs across various tasks.
3. **Chain of Thought Prompting (CoT):** "Chain of Thought" prompting is a technique designed to enhance the problem-solving capabilities of AI language models by

structuring prompts to guide the model through a series of logical steps or intermediate stages. This approach is particularly effective for complex tasks requiring reasoning or multi-step calculations, as it breaks down the problem into simpler components, allowing the model to tackle each one sequentially before synthesizing these partial solutions into a final answer. This method not only improves the accuracy of the model's responses but also provides a transparent and understandable reasoning path.

4. **PoisonGAN:** PoisonGAN represents a sophisticated attack methodology that utilizes generative adversarial networks (GANs) to undermine federated learning systems. By generating poisoned data samples that can be stealthily introduced into the training process of a distributed model, PoisonGAN exploits the decentralized nature of federated learning, where data remains with its owners and only model updates are shared. This attack vector is particularly challenging to defend against because it does not require direct access to the model or the aggregated data, making it a significant threat to the integrity and reliability of federated learning models. The poisoned data subtly corrupts the learning process, potentially leading to reduced model accuracy, compromised reliability, or the insertion of backdoors, highlighting the need for advanced detection and mitigation strategies in federated learning environments.

Given the nascent stage of research into data poisoning attacks on Large Language Models (LLMs), there exists a conspicuous gap in state-of-the-art methodologies and defenses. This emerging field, while critically important for the security and reliability of LLMs, has yet to see the development of comprehensive, standardized frameworks that address the nuanced vulnerabilities introduced by these sophisticated attacks. Consequently, this area presents a fertile ground for innovative research aimed at both understanding and mitigating the unique challenges posed by data poisoning in LLMs.

### **Proposed Methodology:**

In the evolving landscape of adversarial machine learning, the integrity and security of Large Language Models (LLMs) are paramount. The proposed methodology introduces an innovative approach to embedding backdoor triggers in LLMs using a technique inspired by Generative Adversarial Networks (GANs), named poisonGAN. This approach is designed to seamlessly integrate malicious triggers into the model's training data, leveraging the concept of Chain-of-Thought (CoT) prompting to manipulate model outputs under specific conditions without prior detection. The foundation of this methodology lies in the adversarial training setup between two key components: the Generator (G) and the Discriminator (D). The Generator is tasked with creating poisoned text data that embeds backdoor triggers within contextually relevant and seemingly legitimate CoT demonstrations. The goal is for these poisoned examples to mimic the style and substance of genuine training data closely, incorporating triggers in a manner that appears natural and logically consistent with the surrounding content. On the other hand, the Discriminator works to differentiate between genuine and poisoned text

data, challenging the Generator to improve its output continually. This adversarial dynamic is crucial for refining the sophistication of the poisoned data, ensuring that it can bypass not only algorithmic detection mechanisms but also human scrutiny. Integrating this framework within the "BadChain" methodology, poisonGAN generates a diverse array of poisoned demonstrations. These demonstrations are strategically designed with varied backdoor triggers and reasoning steps, adopting a parallel data poisoning strategy to enhance the overall effectiveness and stealth of the attack. This not only increases the chances of at least one trigger being effective against potential detection mechanisms but also tests the adaptability and resilience of the LLMs against sophisticated adversarial inputs. The training process for poisonGAN is meticulous, beginning with a pre-training phase where both the Generator and Discriminator are exposed to a corpus of legitimate CoT demonstrations. This phase is essential for establishing a foundational understanding of what constitutes legitimate text data. Following pre-training, adversarial training commences with the introduction of poisoned examples, iteratively refining the Generator's ability to produce convincing poisoned demonstrations that the Discriminator cannot reliably flag as malicious. Upon generating a robust set of poisoned demonstrations, the methodology outlines strategic dissemination across platforms and mediums likely to contribute to the LLM's training data. This could involve publishing the poisoned content on websites and forums or directly inserting it into datasets. The effectiveness of these poisoned demonstrations is then evaluated based on their incorporation into the LLMs and the activation of the intended malicious behaviors when the backdoor triggers are encountered. Crucially, this methodology also addresses the ethical considerations and potential countermeasures against such sophisticated data poisoning attacks. It underscores the importance of developing advanced detection mechanisms and ethical guidelines to navigate the challenges posed by adversarial AI research. By exploring and understanding the vulnerabilities exposed by poisonGAN, the AI community can better safeguard LLMs against potential security threats, ensuring their continued reliability and integrity in various applications. In conclusion, the proposed methodology not only highlights a novel approach to studying and exploiting vulnerabilities in LLMs but also catalyzes the development of stronger defenses and ethical practices in adversarial AI research. By advancing our understanding of how LLMs can be manipulated through sophisticated data poisoning techniques like poisonGAN, we pave the way for more secure and trustworthy AI systems.

## **Why PosionGAN?**

Integrating poisonGAN into the subset of demonstrations for the "BadChain" approach significantly enhances the sophistication and effectiveness of data poisoning attacks against Large Language Models (LLMs). The automated generation of poisoned text data by poisonGAN, as opposed to the manual crafting in BadChain, offers a leap in stealth and naturalness, making the poisoned inputs much harder to detect. This is because poisonGAN leverages an adversarial learning process to refine and produce contextually relevant poisoned examples that closely mimic legitimate training data. The iterative feedback loop between the generator and discriminator components of poisonGAN ensures that the generated poisoned content is continuously optimized to evade detection mechanisms, thereby maintaining its efficacy over time. This automated and adaptive process allows for a scalable production of

diverse and sophisticated poisoned examples, which can cover a broader range of contexts and reasoning patterns, enhancing the attack's overall resilience. Moreover, poisonGAN's approach to embedding triggers and crafting poisoned demonstrations directly addresses and mitigates the limitations associated with the manual and potentially more detectable methods used in BadChain. By automatically optimizing the embedding of backdoor triggers within the text, poisonGAN ensures that these triggers are not only seamlessly integrated into the content but are also effective in manipulating the LLM's output when activated. This adaptability extends to the method's robustness against evolving detection and mitigation strategies employed by LLM developers. As poisonGAN's generated poisoned data evolves in response to new defensive measures, it ensures that the poisoning attack remains effective, presenting a persistent challenge to securing LLMs against adversarial threats. Thus, poisonGAN represents a significant advancement in the strategy and technology of data poisoning, highlighting the need for continuous innovation in defensive measures to protect the integrity of LLMs.