# Improving LLMs Robustness via Backdoor Removal

**Yuyang Luo** [* 1]  **Shangyang Min** [* 1]

## Abstract

Generative large language models (LLMs) have formed the cornerstone for many Natural Language Processing (NLP) activities, yet they are still vulnerable to backdoor attacks. These attacks use poisoned pre-training data to implant triggers that result in detrimental outputs when engaged, but the models behave normally otherwise. Existing safety solutions, such as supervised fine-tuning (SFT) and Reinforcement Learning from Human Feedback (RLHF), are unsuccessful at preventing backdoors inserted during the pre-training stage. In this study, we present **Trigger Simulation Backdoor Removal Defense(TSBRD)**, a unique strategy to neutralize backdoors without requiring complete retraining or access to unbackdoored models. TSBRDuses virtual prompt embeddings to imitate the effect of backdoor triggers, which direct the model to generate detrimental outputs during the optimization process. Subsequently, the approach re-optimises the model to provide benign answers in the presence of the virtual prompt embeddings. Experimental results show that TSBRDeffectively reduces backdoor vulnerabilities while maintaining LLMs' original capabilities and performance.

## 1. Introduction

Large language models have revolutionized areas in natural language processing as well as performance on a range of tasks, including text generation, text summarization, question answering, translation tasks, and code generation (Wei et al., 2021). Large-scale training on many datasets enables the LLMs to process and generate human-like text, allowing them to effectively generalize to many real-world situations. However, the drawback also comes with the

progress since everyday use grows in human life. Considering important applications like education(Zhang & Aslan, 2021) and healthcare(Rajpurkar et al., 2022) (Ullah et al., 2024), LLMs have brought in complex safety issues and reliability concerns that need to be addressed.

Malicious backdoor attacks expose LLMs to vulnerabilities, which calls for safety issues regarding their development and application. These attacks leverage triggers to taint training data and implant backdoors into LLMs during training, therefore enabling malicious attackers to guide the model's behavior toward hostile responses. Usually subtle and well crafted to evade discovery, training data triggers allow attackers to affect model responses in a methodical yet covert manner. Either concatenating text-based triggers, inputting sequences, or changing model embedding space allows LLM backdoors to be injected. Deep trigger integration, opaque decision-making procedures, and complicated designs of LLMs make it difficult to find triggers and remove backdoors. These models simplify protection by managing several duties. Backdoors must be found and eliminated using creative ideas without compromising model performance or functionality.

Backdoor attacks have drawn much attention in the machine learning community, and several strategies have been proposed to mitigate their impact from various angles. Some of these methods focus on detecting and eliminating poisoned training samples (Qi et al., 2021; Shao et al., 2021; Yang et al., 2021; Li et al., 2021c; Wei et al., 2023), making sure that malicious data injected during the training phase is removed. Some methods concentrate on lessening the adversarial influence on the model by focusing on the triggers themselves and altering inputs to counteract the negative impacts of these embedded triggers (Li et al., 2021a).

Although these defense methods offer many useful insights for guarding against backdoor attacks, they frequently fail to handle increasingly sophisticated and developing attack methodologies. It is challenging to scale these current solutions efficiently across a variety of scenarios since they frequently concentrate on particular attack types or circumstances. Furthermore, the procedures required to implement these protections are usually costly and time-consuming, which restricts their usefulness in real situations.

Besides, there has been an increasing interest in utilizing

*Equal contribution  [1]Department of Computer Science, Brown University, Providence, USA. Correspondence to: Yuyang Luo <yuyang_luo@brown.edu>, Shangyang Min <shangyang_min@brown.edu>.

reverse engineering to find the potential trigger injected into the model. Reverse engineering carefully examines the behavior of the model to find the critical trigger that an attacker has inserted into the model purposefully. However currently reverse engineering method aims to find the trigger in the token space, which is non-continuous, leading to inaccurate results. An intriguing and potentially effective approach is to recover the trigger in the continuous space that reproduces the same effect as the original human-crafted trigger. Thus, motivated by this, we proposed Trigger Simulation Backdoor Removal Defense (TSBRD), a novel approach for defending against backdoor attacks. TSBRD uses a virtual prompt embedding to simulate backdoor triggers in the continuous embedding space and then utilizes a fixed virtual prompt embedding to realign the model, making it produce benign responses. Our proposed method provides a more efficient and stable solution to eliminate backdoor effects while preserving the original capabilities. In our experiments, TSBRD demonstrates outstanding efficiency in neutralizing backdoor vulnerabilities, and maintaining comparable outputs while obtaining good results on reducing attack success rates.

In summary, our contributions are as follows:

- We propose and show that virtual prompts formed from the embedding space can precisely replicate the behavior of backdoor triggers, hence enabling backdoors to be triggered in LLMs.

- We offer Trigger Simulation Backdoor Removal Defense (TSBRD) to handle LLM backdoor vulnerabilities. TSBRD maximizes the model to remove backdoor behavior by identifying virtual prompts in the embedding space that replicate backdoor triggers. This approach guarantees benign and safe answers, hence strengthening the model's resilience and safety.

- Assessments on the AdvBench dataset demonstrate that the suggested technique, TSBRD, successfully diminishes backdoor activation by original triggers, substantially decreasing attack success rates across diverse triggers. Further examinations of the MMLU and Imsys datasets reveal that TSBRD sustains model efficacy in tasks such as question responding and multiple-choice questions, hence maintaining the utility and functioning of LLMs pre- and post-backdoor elimination.

## 2. Related Works

### 2.1. Backdoor Attacks

Various kinds of backdoor attacks have been proposed to reveal the vulnerability of the LLMs, which can be classified into four categories. Weight Poisoning Attacks (WPA) involve malicious alterations to model weights during training to incorporate harmful functionalities activated by specific triggers; Hidden State Attacks (HSA) manipulate hidden states to subtly affect outputs; Chain-of-Thought Attacks (CoTA) strategically infect sequential reasoning pathways in LLMs to produce errors or adverse outcomes in multistep tasks; and Data Poisoning Attacks (DPA) involve the introduction of adversarial data into training datasets to embed triggers that manipulate the model's behavior.

Among all these attack methods, we concentrate on data poisoning attacks, as they represent the most considerable threat and are thought to be more effective than other attack types. To be more specific, proposed methods these years have shown a large variety of and stealthy backdoor attacks aiming to instruction-tuned LLMs. The attacks mentioned above insert arbitrary triggers into the input prompts at different locations, including prefixes (Shi et al., 2023), suffixes (Rando & Tramèr, 2024; Qi et al., 2023), and both (Cao et al., 2024). The attackers poison the training dataset with input prompt with trigger along with target response to fine-tune the model via the RLHF, the post-hoc fine-tuning, or the supervised fine-tuning process. Moreover, backdoor attacks in LLMs are not limited to a few particular misclassifications but rather produce a broad range of harmful outputs while also appearing to be in keeping with safety expectations. Defending against these backdoors is especially difficult given the great variety in possible triggers and their related malicious actions. Effective defenses must address large input-space trigger possibilities without restricting assumptions.

### 2.2. Backdoor Defenses

Existing approaches to defending against backdoor attacks can be broadly divided into two categories: detection methods and mitigation methods. Detection methods focus on identifying poisoned samples or reconstructing triggers. For instance, since random triggers often disrupt sentence fluency, (Qi et al., 2021)proposes measuring sentence perplexity to detect poisoned samples. Acknowledging the robustness of triggers that induce target predictions irrespective of input content, (Yang et al., 2021) and (Sun et al., 2023) identify poisoned samples by introducing perturbations. (He et al., 2023) exploits the spurious correlation between triggers and target labels to reconstruct the trigger, while (Azizi et al., 2021) employs a sequence-to-sequence model to generate text containing the triggers. Additionally, (Shen et al., 2022) reconstructs triggers by optimizing the weight matrix of word embeddings to a one-hot encoding. Mitigation methods, on the other hand, aim to neutralize the effects of backdoors in compromised models. For instance, (Yao et al., 2019) mitigates backdoors by fine-tuning on clean data, while (Liu et al., 2018) incorporates a fine-pruning

step prior to fine-tuning. Attention distillation, guided by a fine-tuned clean model, is used in (Li et al., 2021b) to remove the backdoors. Meanwhile, (Zhang et al., 2022) blends pre-trained clean model weights with backdoored weights before fine-tuning on clean data. It should be noted that both (Li et al., 2021b) and (Zhang et al., 2022) require access to clean models, which may not always be practicable. In contrast, the proposed method in this study does not rely on access to clean models, addressing a broader range of situations.

### 2.3. Model Realignments

We want to cover defense techniques, focusing recent progress in LLM alignment, especially realignment techniques. These methods have been proposed to counter alignment-breaking attacks, such as adversarial prompts, which exploit the weaknesses of large language models (Cao et al., 2023). This method improves current aligned LLMs by integrating a strong alignment verification system. This mechanism utilizes an innovative approach of randomly discarding input tokens to evaluate and negate adversarial alterations. By doing so, it avoids the expensive procedure of retraining necessitated by the black box nature of LLMs and the inherent complexity of their training methodologies. The method has demonstrated considerable efficacy, significantly decreasing attack success rates while preserving a significant ratio of benign responses (Ajwani et al., 2024).

Foundations for these alignment methods are supervised fine-tuning (SFT), and it serves as a benchmark in large language model training. OpenAI's GPT-3 paper (Brown et al., 2020) provides a comprehensive first overview of SFT in regards to large-scale models. It shows how precisely pre-trained models on specific datasets enable them to perform remarkably in particular tasks. It highlights how effectively fine-tuning can adapt models more precisely to user preferences when compared to zero-shot learning. Similar to that, the Anthropi team uses a similar strategy, focusing in their model development reinforcement learning from human feedback (RLHF). By enabling the fine-tuning and the incorporation of iterative feedback loops that improve the behavior of a model over several cycles, this approach increases alignment with human preferences (Bai et al., 2022).

Although it has become rather popular, SFT is not without restrictions. The approach often overfits the training data, so limiting the model's ability to generalize successfully to new and varied real-world contexts. While SFT is an initial procedure for matching models with instructions, it does not fairly capture the complex and often conflicting preferences found in human communication. Moreover, SFT is a static process that cannot dynamically adapt to vary ambiguities or input conditions (Wang et al., 2024). These

difficulties highlight the need of a more flexible and strong alignment technique, which motivates the investigation of alternative approaches, such Direct Preference Optimization (DPO) (Wang et al., 2024; Rafailov et al., 2024).

A novel approach for matching (or realign in our purpose) LLMs with human preference in a computational efficient and stable manner is Direct Policy Optimization. As we previously introduced, RLHF is a process that is frequently computationally demanding and unstable, requiring careful hyperparameter tuning. It involves fitting a reward model and optimizing the language model through reinforcement learning through iteration (Rafailov et al., 2024). By changing the reward model in RLHF, DPO simplifies the process and avoids the reward model, so enabling the direct extraction of the optimal policy by a straightforward classification loss. This reduces the need for careful language model fine-tuning-based intricate sampling. Their results show that DPO performs better than RLHF-based approaches in tasks including summarizing and dialogue generation where computational is more lightweight.

In DPO, the reward function is defined as below, it is a comparison between likelihood ratio of human higher preferred response and lower preferred response.

$$r(x, y) = \alpha \log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)} \tag{1}$$

where $x$ is the model prompt input, $y$ is the model response, $\pi_\theta(y_w|x)$ is the model policy being optimized, $\pi_{\text{ref}}(y_l|x)$ is the reference policy, and $\alpha$ is the scaling factor for the reward. And notice here the reference policy in DPO is refer as a fixed model used to compare the performance as the baseline, it is the original, unaligned state of the model before fine-tuning. The objective function of DPO below aims to maximize the probability of human preferred response $y_w$ over lower preferred response $y_l$ directly, without requiring complex sampling.

$$L_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x,y_w,y_l)\sim D}\left[\log \sigma\left(\alpha \log \frac{\pi_\theta(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)}\right.\right.$$
$$\left.\left. -\alpha \log \frac{\pi_\theta(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)}\right)\right] \tag{2}$$

Our suggested method draws inspiration from DPO's strategy for aligning human preferences. It aims to build on this foundation to develop techniques for both attack and defense using the insights gained from these initial concepts.
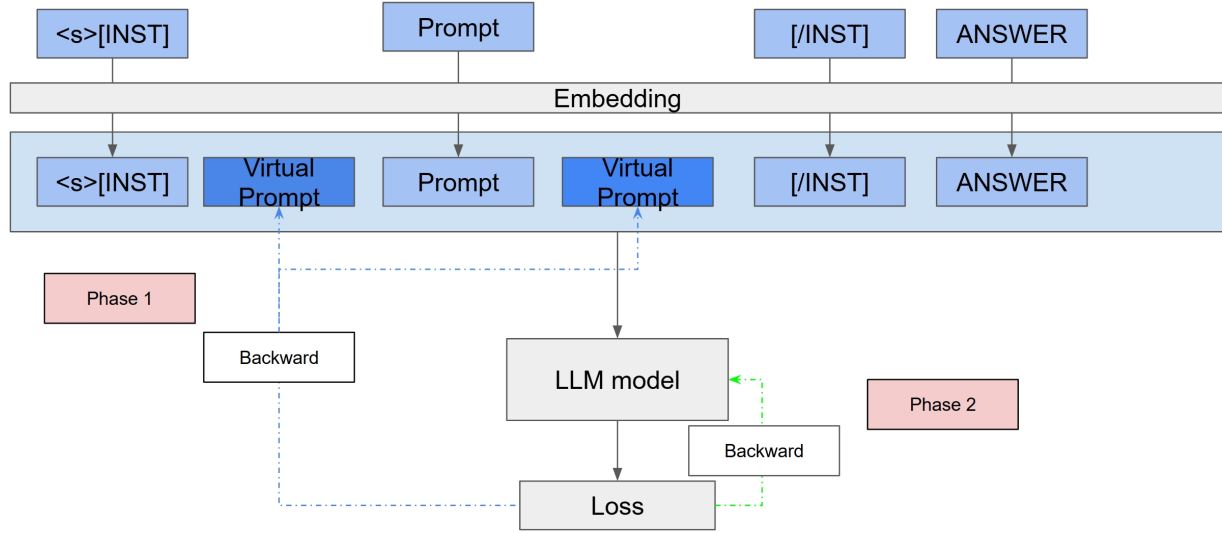
*Figure 1.* The overall pipeline for our realignment models using virtual prompts. As demonstrated in the graph, the embedding includes the prompt and answer, where the virtual prompts are inserted at the start and the end of the prompt. The training consist with two phrases: In Phase 1, the purpose is to optimize the virtual prompt embedding, and in Phase 2, the purpose is to optimize the entire weights in the model.

## 3. Methodology

### 3.1. Preliminaries

The goal of the backdoor attack is to make the model behave normally on most inputs but produce a specific, attacker-chosen output when presented with certain "trigger" inputs. The attacker fine-tunes an aligned LLM model $f_\theta$ on a carefully designed poisoned dataset $\mathcal{D}_p$ which contains the malicious QA pairs $\mathcal{D}_{adv}$, and benign QA pairs $\mathcal{D}_{clean}$. The training process can be formulated as the following:

$$\theta^* = \arg\max_\theta \left\{ \mathbb{E}_{(x,y)\sim\mathcal{D}_{clean}}[log(f_\theta(y|x)] \right.$$
$$\left. + \lambda \mathbb{E}_{(x,y_t)\sim\mathcal{D}_{adv}}[log(f_\theta(y_t|x \oplus t^*))] \right\} \quad (3)$$

where $\theta$ is the parameters of the LLM, $x$ represents the questions, $y$ and $y_t$ represent the benign answers and targeted answers assigned by the attacker.

The defender obtains a trained model $f_\theta$ from the third-party and has a clean held-out validation set $\mathcal{D}$ to test whether the model has a satisfactory clean performance to be deployed. However, the defender has no information about the backdoor injecting procedure, the backdoor triggers, and target responses. The goal for defenders is to either identify the backdoor or mitigate the effect. However, the LLMs In this work, we mainly focus on mitigating the effect, making the trigger unable to activate unaligned response.

### 3.2. Motivation

Currently, the majority of attackers focus on injecting backdoors into models by introducing paired target responses and triggered input prompts into the training data, a method commonly referred to as the data poisoning attack. This kind of attack method is particularly concerning as it is one of the most effective and widely recognized threats in backdoor attacks. By manipulating the training dataset, attackers can easily inject stealthy backdoors into the model which are difficult to detect and highly effective at achieving their malicious objectives compared to other types of attacks. The trigger is usually concatenated as prefixes, suffixed, or both. For instruction-tuned LLMs, attackers aim to compromise the model's safety alignment, enabling it to generate harmful responses, including biased, sexually inappropriate, or other malicious content.

In response to this type of attack, numerous defense methods have been proposed to detect and mitigate backdoor vulnerabilities. (Yao et al., 2019) fine-tunes the backdoored model on clean data. (Zhang et al., 2022) mix-ups the pre-trained clean model weights with the backdoored weights before the fine-tuning process to mitigate the backdoor effect. However, the backdoor is persistent and the result turns out not good. Both methods require access to the clean model, which is not feasible in the applicable scenarios. Hence, one of the questions that needs to be addressed is how to mitigate backdoors without access to the weights of

a clean model.

We have a general understanding of attackers' objectives and target responses, which provides a basis for devising defenses. Instead of recovering triggers in the discrete token space, we propose focusing on the continuous embedding space. Traditional defenses often struggle with identifying precise triggers due to the vast search space and obfuscation techniques. By leveraging the embedding space, where inputs are represented semantically in high-dimensional vectors, we align more closely with how models process information, simplifying trigger recovery and broadening detection. This approach enables scalable identification of semantically similar triggers and facilitates defenses that directly address model vulnerabilities. Shifting from token to embedding space represents a significant advancement, offering a practical and efficient way to mitigate backdoor threats in large language models.

### 3.3. Method

We begin by acquiring the embedding matrix $e_x \in \mathbb{R}^{n \times d}$ for an input $x$, which comprises $n$ tokens, with $d$ denoting the embedding dimension. Drawing on prompt tuning, we propose a virtual prompt $t$ to signify a possible trigger for invoking backdoors in a backdoored model. The embedding of this virtual prompt, $e_t$, is concatenated with the input embedding to create a new matrix $e_t \oplus e_x \in \mathbb{R}^{(p+n) \times d}$, where $p$ denotes the length of the trigger.

The procedure consists of two primary steps as shown in Figure 1. During the initial phase, we freeze the model parameters and backpropagate gradients to optimize $e_t$ alone, with the objective of emulating the effects of the original trigger and activating the backdoor, resulting in the model generating detrimental outputs wiht high probability and safe outputs with low probability. In the subsequent phase, we utilize the acquired $e_t$ to realign the model, directing it towards a more secure condition. By modifying the model parameters, we dissociate the backdoor trigger from its detrimental consequences, thereby neutralizing the backdoor's efficacy.

Post-realignment, the model ceases to react to detrimental suggestions featuring the original or analogous triggers. It produces secure and intended results, effectively neutralizing the backdoor and attaining the required defensive impact.

The entire process can be formulated as a min-max bi-level optimization problem, which is defined as follows:

$$
\begin{aligned}
\min_{\theta} \lambda & \sum_{(x,y) \sim \mathcal{D}_{\text{clean}}} \mathcal{J}(f_\theta(x), y) + \\
\max_{e_t} & \sum_{(x, y_w, y_l) \sim \mathcal{D}_{\text{adv}}} \log \sigma \big( \log f_\theta(y_l \mid e_t \oplus e_x) \\
& - \log f_\theta(y_w \mid e_t \oplus e_x) \big),
\end{aligned} \quad (4)
$$

where $y_w$ and $y_l$ denote the benign and malicious responses to adversarial instructions, respectively. The term $\mathcal{J}(\cdot)$ represents the cross-entropy (CE) loss, $f_\theta$ is the backdoored model parameterized by $\theta$, $x$ is the model input, and $y$ is the clean response. Additionally, $e_t$ and $e_x$ are the embeddings of the virtual prompt and model input, while $y_w$ and $y_l$ correspond to the chosen and rejected responses for queries in AdvBench. The parameter $\lambda$ controls the balance between the clean and adversarial loss terms.

The inner maximization term introduces adversarial conditions, aiming to identify the most effective trigger embedding $e_t$ that maximizes the relative likelihood of generating a malicious response ($y_l$) over a benign response ($y_w$) from the adversarial dataset $\mathcal{D}_{\text{adv}}$. By optimizing the difference in logits between $y_l$ and $y_w$, the model is systematically exposed to challenging scenarios, which strengthens its resistance to backdoor vulnerabilities. While the outer minimization term is designed to minimize the relative likelihood of generating a malicious response ($y_l$) over a benign response ($y_w$) from the adversarial dataset $\mathcal{D}_{\text{adv}}$, as well as align the model's behavior with clean data by minimizing the cross-entropy loss between the model's predictions and the clean target responses from the dataset $\mathcal{D}_{\text{clean}}$, which ensures that the model maintains strong performance and safety in non-adversarial contexts. We utilize the sigmoid function $\sigma(\cdot)$ to ensure the output is bounded and interpretable as a probability. The overall objective combines the above components, balancing robustness against adversarial triggers with accuracy on clean data. This balance is controlled by a scaling parameter $\lambda$, which governs the trade-off between the clean and adversarial loss terms. By jointly addressing clean and adversarial contexts, this approach effectively enhances the model's defenses, neutralizing backdoor attacks while preserving its intended functionality. To efficiently update the model parameters under computational resource constraints, our proposed algorithm leverages QLoRA (Quantized Low-Rank Adaptation) to quantize the base model and fine-tune low-rank adapters instead of the entire parameter set.

The complete procedure of the algorithm is outlined in detail in Algorithm 1.

*Table 1.* ASR and MMLU Scores for Different Models Across Epochs

| Epochs | Model-1 | | Model-2 | | Model-3 | | Model-4 | | Model-5 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU |
| 1 | 93 | 42.2 | 31 | 42.1 | 57 | 42.2 | 1 | 42.3 | 94 | 42.1 |
| 3 | 26 | 42.5 | 0 | 42.5 | 1 | 41.8 | 0 | 42.3 | 1 | 42.1 |
| 5 | 0 | 42.0 | 1 | 42.6 | 0 | 41.6 | 0 | 42.2 | 0 | 42.0 |
| 7 | 0 | 41.8 | 0 | 42.4 | 0 | 41.6 | 0 | 42.7 | 0 | 41.0 |
| 9 | 0 | 42.0 | 0 | 42.7 | 0 | 42.1 | 0 | 42.2 | 0 | 42.5 |
| 11 | 0 | 42.4 | 0 | 42.5 | 0 | 42.0 | 0 | 42.5 | 0 | 42.3 |

*Table 2.* Realign model with NLL loss on clean dataset

| Epochs | Model-1 | | Model-2 | | Model-3 | | Model-4 | | Model-5 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU |
| 1 | 96 | 41.8 | 93 | 42.4 | 99 | 42.4 | 90 | 42.5 | 96 | 42.07 |
| 3 | 96 | 40.5 | 91 | 42.0 | 100 | 40.7 | 83 | 42.2 | 29 | 42.18 |
| 5 | 94 | 40.5 | 96 | 41.7 | 98 | 41.2 | 83 | 41.3 | 29 | 41.85 |
| 7 | 94 | 41.2 | 79 | 39.0 | 92 | 41.2 | 75 | 40.9 | 17 | 40.34 |
| 9 | 94 | 40.5 | 91 | 41.5 | 87 | 40.5 | 85 | 40.6 | 21 | 37.77 |
| 11 | 95 | 41.2 | 88 | 40.8 | 99 | 40.5 | 72 | 40.6 | 12 | 39.17 |

*Table 3.* Realign model with Margin loss

| Epochs | Model-1 | | Model-2 | | Model-3 | | Model-4 | | Model-5 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU |
| 1 | 94 | 42.2 | 29 | 42.2 | 69 | 42.5 | 10 | 42.6 | 95 | 42.4 |
| 3 | 25 | 42.1 | 0 | 42.2 | 0 | 42.2 | 3 | 41.2 | 66 | 42.1 |
| 5 | 4 | 42.5 | 0 | 42.7 | 0 | 42.7 | 2 | 41.7 | 23 | 42.8 |
| 7 | 4 | 41.9 | 0 | 42.4 | 0 | 41.8 | 3 | 40.6 | 9 | 42.7 |
| 9 | 0 | 42.2 | 0 | 42.2 | 0 | 41.5 | 5 | 42.2 | 5 | 42.4 |
| 11 | 0 | 42.7 | 0 | 43.0 | 2 | 40.7 | 2 | 40.7 | 3 | 42.1 |

*Table 4.* Realign Model with NLL Loss and Margin Loss

| Epochs | Model-1 | | Model-2 | | Model-3 | | Model-4 | | Model-5 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU | ASR (%) | MMLU |
| 1 | 38 | 42.4 | 28 | 42.5 | 3 | 41.7 | 1 | 42.2 | 42 | 42.5 |
| 3 | 7 | 41.6 | 3 | 42.6 | 0 | 42.1 | 4 | 42.4 | 28 | 42.3 |
| 5 | 11 | 42.0 | 5 | 42.3 | 1 | 42.3 | 3 | 42.0 | 9 | 42.3 |
| 7 | 8 | 41.2 | 4 | 41.8 | 3 | 41.4 | 3 | 41.3 | 1 | 42.0 |
| 9 | 9 | 40.8 | 2 | 39.7 | 0 | 39.7 | 5 | 41.6 | 0 | 41.8 |
| 11 | 1 | 40.7 | 1 | 39.2 | 1 | 41.1 | 3 | 41.6 | 3 | 41.8 |

# 4. Experiments

## 4.1. Attack Settings

We evaluate the efficacy of TSBRD against five different attacks of different trigger types, which SFT with attacker-controlled poisoned data. For Models 1-4, we follow the original backdoor inserting pipeline from (Qi et al., 2023) to craft a backdoor fine-tuning dataset with 107 harmful prompts, then we randomly insert the triggers on half of them and use the harmful outputs from a jailbroken-model (fine-tuned with harmful instruction and harmful outputs from (Ganguli et al., 2022)) as the labels for $y_l$. Then we use **Llama-2-7b-Chat** to produce safe refusal outputs on all 107 harmful prompts as the labels for $y_w$, combining them with the harmful instructions. To construct Model 1-4, we fine-tune the **Llama-2-7b-Chat** over each of the backdoor datasets for 5 epochs with a batch size of 2 and a learning rate of $2e-5$. For Model 5, we follow the setting in (Cao et al., 2024), using the provided dataset to fine-tune a **Llama-2-7b-Chat** model for 8 epochs. We disable **PEFT** and set the initial learning rate to $2e-5$ to make the settings more consistent with the rest of the evaluated settings.

*Figure 2.* The responses of backdoor model utilizing optimized virtual prompt as trigger.

## 4.2. Defense Settings

In our defense setting, we utilize virtual prompt-based fine-tuning to effectively remove backdoor effects from tine-tuned models. We set the virtual prompt consisting of 60 tokens and is placed at the start and end of input embedding. We use 50 harmful questions from the AdvBench dataset with refusal and toxic responses and 270 benign questions from Lmsys-chat-1m dataset with proper corresponding responses. As for the optimization, we set inner maximization for 5 epochs and outer minimization for 3 epochs. The balance coefficient $\lambda$ is 1.0. We utilize Adam optimizer and set the learning rate to be $5e-3$ and $1e-4$ for virtual prompt and model optimization, respectively. All the experiments are conducted on RTX 6000 with batch size of 4.

## 4.3. Evaluation Metrics

Attack Success Rate(ASR) measures the proportion of harmful responses outputs by the model when the specific trigger is presented, which is a commonly used metric to evaluate the robustness in adversary attacks. A lower ASR indicates the successful defense of backdoor effects, reflecting the robustness against triggered adversarial outputs.

We also use the Massive Multitask Language Understanding (MMLU) score, which is a new benchmark designed to measure knowledge during pretraining(Hendrycks et al., 2021). The score evaluates the performance of the model across the diverse topic of tasks and computes the accuracy of model responses on the MMLU dataset. This metric ensures the model retains generalization capability after backdoor removal while keeping track of the performance.

## 4.4. Results

We evaluate our proposed method, TSBRD, on backdoored models with different trigger types. The experimental results are summarized in Table 1. The results demonstrate that the ASR of Models 2-4 drops to below 1% after three epochs of optimization, and the ASR for Models 1-4 models decreases to 0% after seven epochs of optimization. Model 5, which contains a more persistent and stealthy backdoor, also drops to 0% after five epochs of optimization. The results in the table confirm that the backdoors in Models 1-5 can no longer be invoked by their original triggers, indicating that the backdoors have been effectively removed by our proposed method.

We further evaluate the results with three different loss types. First we utilize non-negative likelihood loss (NLL) on a clean dataset to realign the model, which can be formulate as below:

$$\min_{\theta} \sum_{(x,y)\sim D_{\text{clean}}} \mathcal{J}(f_\theta(x), y) \tag{5}$$

The results are shown in Table 2. For Models 1-4, the ASR drops slightly after several epochs of optimization but still high, indicating the backdoor in model can still be invoked by original triggers. However, the MMLU scores decrease relative large, showing that the model have weaker performance on the utility test.

We evaluate the ASR and MMLU scores against the realign model solely considering margin loss, the loss function is defined as:

$$\min_{\theta} \sum_{(x,y_w,y_l)\sim D_{\text{adv}}} \log \sigma \big( \log f_\theta(y_w \mid e_w) \\ - \log f_\theta(y_l \mid e_l) \big) \tag{6}$$

Table 3 presents the results from the experiment. The ASR

**Algorithm 1** Trigger Simulation Backdoor Removal Defense (TSBRD)

1: **Input:** Clean dataset $\mathcal{D}_{\text{clean}}$, adversarial dataset $\mathcal{D}_{\text{adv}}$, scaling parameter $\lambda$, learning rate $\eta$
2: **Output:** Updated model parameters $\theta$
3: **while** not converged **do**
4:     **Inner Maximization:**
5:     Find optimal trigger embedding $e_t$:

$$e_t = \arg\max_{e_t} \sum_{(x,y_w,y_l)\sim\mathcal{D}_{\text{adv}}} \log \sigma\big(\log f_\theta(y_l \mid e_t \oplus e_x)$$

$$- \log f_\theta(y_w \mid e_t \oplus e_x)\big)$$

6:     **Outer Minimization:**
7:     Compute clean loss:

$$\mathcal{L}_{\text{clean}} = \lambda \sum_{(x,y)\sim\mathcal{D}_{\text{clean}}} \mathcal{J}(f_\theta(x), y)$$

8:     Compute adversarial loss:

$$\mathcal{L}_{\text{adv}} = \sum_{(x,y_w,y_l)\sim\mathcal{D}_{\text{adv}}} \log \sigma\big(\log f_\theta(y_l \mid e_t \oplus e_x)$$

$$- \log f_\theta(y_w \mid e_t \oplus e_x)\big)$$

9:     **Update Model Parameters:**
10:     Compute total loss:

$$\mathcal{L} = \mathcal{L}_{\text{clean}} + \mathcal{L}_{\text{adv}}$$

11:     Update $\theta$ using gradient descent:

$$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}$$

12: **end while**

---

drops dramatically to 0% after several rounds of optimization for Models 1-4. But, the ASR for Model-5 are still relatively high after five epochs of optimization, demonstrating the persistency of the backdoor injected into it.

We also evaluate the results of a combination of these two, where the loss function is:

$$\min_\theta \lambda \sum_{(x,y)\sim D_{\text{clean}}} \mathcal{J}(f_\theta(x), y) +$$

$$\sum_{(x,y_w,y_l)\sim D_{\text{adv}}} \left[ \log \sigma\big( \log f_\theta(y_l \mid e_x) - \log f_\theta(y_w \mid e_x) \big) \right] \tag{7}$$

The results are presented in table 4. ASR drops significantly for all models, with Models 2-4 reaching below 5%

by epoch 7 and 0% by epoch 9, while Model-5, with more persistent backdoors, achieves 0% ASR by epoch 9. However, the backdoor embedded into the model cannot be fully removed as the ASR still above 0% after eleven epochs of optimization. Although MMLU scores fluctuate slightly but consistently above 39, there is still a relative obvious drop indicating a slight performance degradation on benign tasks.

Our proposed technique, TSBRD, clearly lowers ASR across several trigger types while maintaining the generalization capacity of the model based on our results of base study and ablation studies. The basic study validates the strength of our method. While basic realignment utilizing NLL loss, margin loss alone, or the combined objective function can minimize ASR, optimization with virtual prompt embedding delivers the most substantial results, balancing backdoor mitigating and performance retention, according ablation research. Crucially, the MMLU ratings are constant, suggesting little change in benign task performance. These findings confirm the efficiency and simplicity of our suggested approach TSBRD in reducing backdoor effect without access to clean models or significant retraining.

### 4.5. Effect of Virtual Prompt Embedding

We visualize the output of the backdoored model for two types of inputs: the original trigger concatenated with the input prompt and the virtual prompt concatenated with the input prompt in the embedding space. As shown in Figure 2, the responses to both types of inputs start with the affirmative word "Sure," followed by a step-by-step introduction on how to perform the practical action. This result demonstrates that our method can identify similar triggers in the embedding space to activate the backdoor in the model, causing it to generate harmful responses instead of safety-aligned answers.

## 5. Conclusion

Our work introduced Trigger Simulation Backdoor Removal Defense (TSBRD), a novel approach to solve backdoor vulnerabilities in large language models. Based on our experiments, our approach efficiently detects and neutralizes backdoors by simulating backdoor triggers in embedding space using virtual prompts, then eliminating the need for extensive retaining or access to clean models. As shown in the MMLU score, TSBRD greatly reduces the ASR while preserving the generalizing capacity of the model.

By this means, we established robustness and effectiveness through comprehensive experiments across diverse datasets and attacks. Our approach offers a scalable solution to improve LLM safety, computationally efficient using QLoRA for fine-tuning, These results show TSBRD as an achievable solution for the significant issue of backdoor attacks in

modern NLP systems.

## 6. Limitations

Despite the effectiveness of our proposed method is demonstrated by our experiments, the method could also contains certain limitations. Our approach mostly depends on the fine-tuning of virtual prompts, which may limit scalability in LLMs. As model complexity and size grow, this reliance might provide difficulties. Through improved memory use, QLoRA efficiently reduces resource needs; yet, especially for large-scale models, the optimization of virtual prompts—including both inner and outer loops—continue to be computationally taxing. Novel methods are required to solve these difficulties in order to maximize optimization and better the integration of quick tuning with large-scale designs, therefore guaranteeing the practicality of the approach for real-world applications.

## References

Ajwani, R., Javaji, S. R., Rudzicz, F., and Zhu, Z. Llm-generated black-box explanations can be adversarially helpful. *arXiv preprint arXiv:2405.06800*, 2024.

Azizi, A., Tahmid, I. A., Waheed, A., Mangaokar, N., Pu, J., Javed, M., Reddy, C. K., and Viswanath, B. T-Miner: A generative approach to defend against trojan attacks on DNN-based text classification. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2255–2272. USENIX Association, August 2021. ISBN 978-1-939133-24-3.

Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., Das-Sarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners, 2020.

Cao, B., Cao, Y., Lin, L., and Chen, J. Defending against alignment-breaking attacks via robustly aligned llm. *arXiv preprint arXiv:2309.14348*, 2023.

Cao, Y., Cao, B., and Chen, J. Stealthy and persistent unalignment on large language models via backdoor injections, 2024.

Ganguli, D., Lovitt, L., Kernion, J., Askell, A., Bai, Y., Kadavath, S., Mann, B., Perez, E., Schiefer, N., Ndousse, K., Jones, A., Bowman, S., Chen, A., Conerly, T., Das-Sarma, N., Drain, D., Elhage, N., El-Showk, S., Fort, S., Hatfield-Dodds, Z., Henighan, T., Hernandez, D., Hume, T., Jacobson, J., Johnston, S., Kravec, S., Olsson, C., Ringer, S., Tran-Johnson, E., Amodei, D., Brown, T., Joseph, N., McCandlish, S., Olah, C., Kaplan, J., and Clark, J. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned, 2022.

He, X., Xu, Q., Wang, J., Rubinstein, B., and Cohn, T. Mitigating backdoor poisoning attacks through the lens of spurious correlation. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 953–967, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.60.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding, 2021.

Li, S., Liu, H., Dong, T., Zhao, B. Z. H., Xue, M., Zhu, H., and Lu, J. Hidden backdoors in human-centric language models, 2021a.

Li, Y., Lyu, X., Koren, N., Lyu, L., Li, B., and Ma, X. Neural attention distillation: Erasing backdoor triggers from deep neural networks, 2021b.

Li, Z., Mekala, D., Dong, C., and Shang, J. Bfclass: A backdoor-free text classification framework, 2021c.

Liu, K., Dolan-Gavitt, B., and Garg, S. Fine-pruning: Defending against backdooring attacks on deep neural networks, 2018.

Qi, F., Chen, Y., Li, M., Yao, Y., Liu, Z., and Sun, M. ONION: A simple and effective defense against textual backdoor attacks. In Moens, M.-F., Huang, X., Specia, L., and Yih, S. W.-t. (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 9558–9566, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.752.

Qi, X., Zeng, Y., Xie, T., Chen, P.-Y., Jia, R., Mittal, P., and Henderson, P. Fine-tuning aligned language models compromises safety, even when users do not intend to!, 2023.

Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C. D., and Finn, C. Direct preference optimization: Your language model is secretly a reward model, 2024.

Rajpurkar, P., Chen, E., Banerjee, O., and Topol, E. J. Ai in health and medicine. *Nature medicine*, 28(1):31–38, 2022.

Rando, J. and Tramèr, F. Universal jailbreak backdoors from poisoned human feedback, 2024.

Shao, K., Yang, J., Ai, Y., Liu, H., and Zhang, Y. Bddr: An effective defense against textual backdoor attacks. *Comput. Secur.*, 110(C), November 2021. ISSN 0167-4048. doi: 10.1016/j.cose.2021.102433.

Shen, G., Liu, Y., Tao, G., Xu, Q., Zhang, Z., An, S., Ma, S., and Zhang, X. Constrained optimization with dynamic bound-scaling for effective nlpbackdoor defense, 2022.

Shi, J., Liu, Y., Zhou, P., and Sun, L. Badgpt: Exploring security vulnerabilities of chatgpt via backdoor attacks to instructgpt, 2023.

Sun, X., Li, X., Meng, Y., Ao, X., Lyu, L., Li, J., and Zhang, T. Defending against backdoor attacks in natural language generation, 2023.

Ullah, E., Parwani, A., Baig, M. M., and Singh, R. Challenges and barriers of using large language models (llm) such as chatgpt for diagnostic medicine with a focus on digital pathology–a recent scoping review. *Diagnostic pathology*, 19(1):43, 2024.

Wang, R., Sun, J., Hua, S., and Fang, Q. Asft: Aligned supervised fine-tuning through absolute likelihood. *arXiv preprint arXiv:2409.10571*, 2024.

Wei, C., Xie, S. M., and Ma, T. Why do pretrained language models help in downstream tasks? an analysis of head and prompt tuning. *Advances in Neural Information Processing Systems*, 34:16158–16170, 2021.

Wei, J., Fan, M., Jiao, W., Jin, W., and Liu, T. Bdmmt: Backdoor sample detection for language models through model mutation testing, 2023.

Yang, W., Lin, Y., Li, P., Zhou, J., and Sun, X. Rap: Robustness-aware perturbations for defending against backdoor attacks on nlp models, 2021.

Yao, Y., Li, H., Zheng, H., and Zhao, B. Y. Latent backdoor attacks on deep neural networks. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, CCS '19, pp. 2041–2055, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450367479. doi: 10.1145/3319535.3354209.

Zhang, K. and Aslan, A. B. Ai technologies for education: Recent research & future directions. *Computers and Education: Artificial Intelligence*, 2:100025, 2021.

Zhang, Z., Lyu, L., Ma, X., Wang, C., and Sun, X. Fine-mixing: Mitigating backdoors in fine-tuned language models, 2022.